

UNIVERSITY OF HAWAII
LIBRARY
ARCHIVE
JAN 7 '60
for

RATIONAL MECHANICS

and

ANALYSIS

Edited by
C. TRUESDELL

Volume 4, Number 1



SPRINGER-VERLAG
BERLIN-GÖTTINGEN-HEIDELBERG
(Postverlagsort Berlin • 4. 11. 1959)

Mechanicam vero duplicem Veteres constituerunt: Rationalem quae per Demonstrationes accurate procedit, & Practicam. Ad practicam spectant Artes omnes Manuales, a quibus utique Mechanica nomen mutuata est. Cum autem Artifices parum accurate operari soleant, fit ut Mechanica omnis a Geometria ita distinguatur, ut quicquid accuratum sit ad Geometriam referatur, quicquid minus accuratum ad Mechanicam. Attamen errores non sunt Artis sed Artificum. Qui minus accurate operatur, imperfectior est Mechanicus, & si quis accuratissime operari posset, hic foret Mechanicus omnium perfectissimus.

NEWTON

La généralité que j'embrasse, au lieu d'éblouir nos lumieres, nous découvrira plutôt les véritables loix de la Nature dans tout leur éclat, & on y trouvera des raisons encore plus fortes, d'en admirer la beauté & la simplicité.

EULER

Ceux qui aiment l'Analyse, verront avec plaisir la Mécanique en devenir une nouvelle branche ...

LAGRANGE

The ARCHIVE FOR RATIONAL MECHANICS AND ANALYSIS nourishes the discipline of mechanics as a deductive, mathematical science in the classical tradition and promotes pure analysis, particularly in contexts of application. Its purpose is to give rapid and full publication to researches of exceptional moment, depth, and permanence.

Each memoir must meet a standard of rigor set by the best work in its field. Contributions must consist largely in original research; on occasion, an expository paper may be invited.

English, French, German, Italian, and Latin are the languages of the Archive. Authors are urged to write clearly and well, avoiding an excessively condensed or crabbed style.

Manuscripts intended for the Archive should be submitted to an appropriate member of the Editorial Board.

The ARCHIVE FOR RATIONAL MECHANICS AND ANALYSIS appears in numbers struck off as the material reaches the press; five numbers constitute a volume. Subscriptions may be entered through any agent. The price is DM 96.—per volume.

Notice is hereby given that for all articles published exclusive rights in all languages and countries rest with Springer-Verlag. Without express permission of Springer-Verlag, no reproduction of any kind is allowed.

For each paper 75 offprints are provided free of charge.

5856-247

Invariant and Complete Stress Functions for General Continua

C. TRUESDELL

Contents

	Page
1. Definition of stress functions. Completeness	1
2. Nature and purpose of this report	2
3. Mechanical situations to which the results may be applied	3
Part I. Flat spaces	
4. General solution for a flat space of n dimensions	6
5. Three-dimensional flat space	7
6. Application to the axially symmetric case	9
7. Application to unsteady motion subject to plane stress	10
8. Unsteady motion in three dimensions	11
9. Status of the problem for flat spaces	12
Part II. Curved spaces	
10. Formal intrinsic solution in an affine space	13
11. The form of the conditions of compatibility in a Riemannian space	15
12. Application to spaces of constant curvature	18
13. Application to the classical theory of membranes	18
14. Status of the problem for curved spaces	20
References	21
Appendix: Bibliography of works on stress functions for linear elasticity and related theories	24

1. Definition of stress functions. Completeness. If a particular rule enables us to find finite algebraic combinations of the derivatives of a set of arbitrary functions, possibly supplemented by quantities associated with the geometry of the space in question, such that these combinations when substituted for the stress tensor satisfy the equations of equilibrium or motion identically in the arbitrary functions, the rule is said to furnish a *solution in terms of stress functions*. Such solutions are of two kinds:

1. Solutions of particular dynamical equations, such as those of linear elasticity, classical fluid mechanics, *etc.*
2. Solutions of the dynamical equations valid for all continuous media in a particular space.

In the above statement, the "arbitrary" functions are restricted to belong to a certain class, such as functions which are four times continuously differentiable, or harmonic, or biharmonic, but the class must not be specified in terms of material properties, nor may it coincide with the class of solutions of the original equations, for otherwise the definition would be vacuous.

By their definition, stress functions provide sufficient conditions for solving the dynamical equations in question. From any such rule, by particularizing the arbitrary functions we may derive any number of special exact solutions. Thus a solution in stress functions is useful for experimentation in search of particular cases, for inverse and semi-inverse methods, *etc.* But this is not enough. We should prove also that the solution so obtained is *general*, or *complete*: Corresponding to any stress field satisfying the given equations there exists at least one suitable choice of the stress functions — in a word, that the rule expresses a necessary as well as a sufficient condition. Only on the basis of a theorem of completeness are we justified in inferring from use of the stress functions any quality of the theory itself. With such a theorem, however, we may, at least in principle, lay aside the differential equations once and for all and enjoy the advantages of explicit, finite formulae.

The second kind of stress functions is obviously of a more embracing nature than the first. However, when a solution of the second kind has been found, it does not, in general, yield a solution in stress functions for any particular material; rather, when a solution of the second kind is substituted into the relevant constitutive equations, there results, usually, a system of differential or integro-differential or functional equations to be satisfied by the stress functions. This corresponds to the very nature of constitutive equations, which single out among all dynamically possible stresses those compatible with a particular kind of material behavior.

The process just mentioned is not often an efficient one for solution in a particular theory. For example, in linear elasticity the stress functions of GALERKIN and PAPKOVICH, which are stress functions of the first kind, furnish a more practical approach to the solution of special problems and the proof of general theorems. But the moment that linear elasticity and closely related theories are left behind, these stress functions fail altogether. Stress functions of the second kind, however, are valid independently of constitutive equations. They should not be looked upon so much as tools to be applied in any special theory but as means of finding *properties common to all theories* and of *comparing the results of different theories*. The rise of various non-linear theories of materials in the past decade suggests that applications of this kind, though presently few, may soon be much demanded¹.

2. Nature and purpose of this report. Here we consider only the second kind of solution, stress functions for general continua², independent of constitutive assumptions. Since the ever growing literature, while giving evidence of broad and continuing interest, includes many rediscoveries, it has seemed worthwhile to put in one place all major known results and methods. The derivations are carried through in general co-ordinates, yielding results in *invariant form*, and the solutions are proved to be *complete* within a stated class.

¹ Stress functions do not furnish the only possibility of such comparisons. In fact, most of the presently known results of this kind have been obtained by SIGNORELLI's theory of stress means; for an exposition, see §§ 220—222 of the forthcoming work by TRUESDELL & TOUPIN, *The classical field theories*, Handbuch der Physik III/1, Springer-Verlag.

² A bibliography of works on stress functions for the classical linear theory of elasticity is subjoined at the end of the report.

The equations to be solved, in the most general case envisioned, are

$$T^{km}_{,m} = 0, \quad T^{km} = T^{mk}, \quad (2.1)$$

where the quantities T^{km} are the contravariant components of an absolute tensor field and where the comma denotes covariant differentiation based upon an affine connection in a real space of n dimensions.

The special character of many results recently published reflects haphazard methods, unnecessary had use been made of either of the two known general approaches:

1. Representation of a solenoidal field as the curl of a vector potential.
2. Use of the principle of virtual work and its converse.

Both methods should be applicable to curved spaces. That the former method has been used only for flat spaces reflects the fact that the general theory of potentials is unknown for curved spaces. In a flat space of n dimensions, the equation $T^{k_1 \dots k_p m}_{,m} = 0$ may be regarded as a system of n^p independent vector equations of the type $V^m_{,m} = 0$, so that the ordinary theorem of the vector potential suffices³. In a curved space, the equations $T^{km}_{,m} = 0$ are not known to be reducible to a system of equations of the form $V^m_{,m} = 0$, so the fact that the general solution of $V^m_{,m} = 0$ is known has not led to solution of our problem. The latter method, in its first stages, applies equally easily to curved spaces. To get fully explicit results by its means, however, the conditions of compatibility in the appropriate space must be exhibited. Calculation of the conditions of compatibility and proof of the theorem of the vector potential are thus seen to be closely related problems, neither of which has been solved for a general curved space. For ordinary curved surfaces, however, enough is known that we can illustrate the power of the latter method by obtaining the most general equilibrated stresses in several classes of curved membranes.

Both methods, in contrast to the more popular method of hit-or-miss trials, yield or imply proofs of completeness.

3. Mechanical situations to which the results may be applied. In an inertial co-ordinate system in a real Euclidean space of three dimensions, the stress tensor t^{km} satisfies CAUCHY's laws of motion

$$\rho \ddot{x}^k = t^{km}_{,m} + \rho f^k, \quad t^{km} = t^{mk}. \quad (3.1)$$

These laws express the local balance of linear momentum and of moment of momentum, respectively, on the supposition that there are neither extrinsic couples nor couples regarded as representing the actions which maintain material continuity. The tensor \mathbf{t} is presumed to include any electromagnetic stresses that may be present. ρ is the density, $\ddot{\mathbf{x}}$ the acceleration, and \mathbf{f} the extrinsic force per unit mass.

CAUCHY's laws (3.1) assume the form (2.1) in various special cases.

Case 1. *Internally free three-dimensional medium in equilibrium.* If $\dot{\mathbf{x}}$, the velocity, is constant and uniform, and if $\mathbf{f} = 0$, (3.1) assumes the form (2.1) with $\mathbf{T} = \mathbf{t}$.

³ Cf. FINZI & PASTORI [1949, 1, Cap. IV, § 7].

Case 2. *Equilibrium of a three-dimensional medium subject to conservative extrinsic force.* If for some scalar $U(\mathbf{x}, t)$ we have

$$\varrho \dot{f}_k = -U_{,k}, \quad (3.2)$$

and again $\dot{\mathbf{x}} = \text{const.}$, (3.1) assumes the form (2.1) if we put $\mathbf{T} = \mathbf{t} - U \mathbf{1}$.

In any case, if the general solution of $t^{k,m}_{,m} = 0$ is known, all that is needed in order to find the general solution of $t^{k,m}_{,m} = -\varrho \dot{f}^k$ is to add a particular integral. In rectangular Cartesian co-ordinates, such an integral $p^{k,m}$ is furnished for convex domains by

$$\|p^{k,m}\| = \left\| \begin{array}{ccc} -\int \varrho f_x dx & 0 & 0 \\ 0 & -\int \varrho f_y dy & 0 \\ 0 & 0 & -\int \varrho f_z dz \end{array} \right\|. \quad (3.3)$$

In curved spaces, it may be more difficult to find a particular integral. For a Riemannian space, the search for a particular integral may be systematized as follows⁴: Let u^k be any contravariant vector satisfying the partial differential system

$$u^{k,m}_{,m} - R^k_m u^m + \varrho (\dot{f}^k - \ddot{x}^k) = 0, \quad (3.3a)$$

where R^k_m is the contracted curvature tensor; then a particular integral $p^{k,m}$ of the system (3.1) is given by

$$p^{k,m} = 2u^{(k,m)} - u^p_p g^{k,m}. \quad (3.3b)$$

The problem of finding a particular integral of the underdetermined tensorial equations is thus reduced to that of finding a particular integral of a determinate vectorial system of Helmholtz type: $\nabla^2 \mathbf{u} = f(\mathbf{x}) \mathbf{u} + \mathbf{g}(\mathbf{x})$.

Case 3. *Steady motion of a three-dimensional medium.* In virtue of EULER'S equation of continuity,

$$\frac{\partial \varrho}{\partial t} + (\varrho \dot{x}^k)_{,k} = 0, \quad (3.4)$$

we may put CAUCHY'S first law (3.1)₁ into the form

$$\frac{\partial (\varrho \dot{x}^k)}{\partial t} = (t^{k,m} - \varrho \dot{x}^k \dot{x}^m)_{,k} + \varrho \dot{f}^k. \quad (3.5)$$

In the case of steady motion subject to no extrinsic force, (3.5) reduces to (2.1)₁ if we put $\mathbf{T} = \mathbf{t} - \varrho \dot{\mathbf{x}} \dot{\mathbf{x}}$, and (2.1)₂ is satisfied in virtue of (3.1)₂.

Cases 2 and 3 may be combined.

Case 4. *Unsteady motion of a three-dimensional medium.* Let indices Γ, Δ run from 1 to 4, and retain k, m as indices with range 1 to 3. Put

$$x^\Gamma \equiv (x^k, t), \quad V^\Gamma \equiv (\dot{x}^k, 1), \quad (3.6)$$

and

$$\|t^{\Gamma\Delta}\| \equiv \left\| \begin{array}{cc} t^{k,m} & 0 \\ 0 & 0 \end{array} \right\|, \quad (3.7)$$

and define a stress-momentum tensor \mathbf{T} as follows:

$$T^{\Gamma\Delta} \equiv t^{\Gamma\Delta} - \varrho V^\Gamma V^\Delta. \quad (3.8)$$

⁴ For a flat space, the method is due to SCHAEFER [1953, 5, § 5].

Thus, e.g., $T^{T4} = -\rho V^T = T^{4T}$, $T^{44} = -\rho$. Suppose that $\mathbf{f} = 0$, and select *rectangular Cartesian space co-ordinates* x^k . Then (3.4) and (3.5) together assume the form (2.1)₁, where the comma now stands for partial differentiation. Thus the laws of motion in an inertial frame assume the form to which (2.1) would reduce *if space-time were flat*. This observation suffices for our later analysis of unsteady motion, which does not presume any geometrical structure of space-time⁵.

Case 5. *Plane motion subject to plane stress.* When $\dot{\mathbf{x}}$ and \mathbf{t} are plane, all the foregoing cases become applicable, but the dimension of the space drops by one. In Case 4, for example, the equations of plane unsteady motion assume the same form, in an inertial frame, as the equations of equilibrium in a three-dimensional flat space.

The equations of equilibrium of a shell having first and second fundamental forms \mathbf{a} and \mathbf{b} , respectively, are of the form⁶

$$\begin{aligned} S^{\delta}_{,\delta} + b_{\gamma\delta} S^{\gamma\delta} + F &= 0, \\ S^{\gamma\delta}_{,\delta} - a^{\gamma\delta} b_{\sigma\delta} S^{\delta} + F^{\delta} &= 0, \\ M^{\delta}_{,\delta} + b_{\gamma\delta} M^{\gamma\delta} + e_{\gamma\delta} S^{\gamma\delta} + L &= 0, \\ M^{\gamma\delta}_{,\delta} - a^{\gamma\sigma} b_{\sigma\delta} M^{\delta} + a^{\gamma\delta} e_{\delta\sigma} S^{\sigma} + L^{\gamma} &= 0, \end{aligned} \quad (3.9)$$

where the comma denotes covariant differentiation based upon the surface metric \mathbf{a} , where the indices have the range 1, 2, and where $e_{\gamma\delta} \equiv +\sqrt{a} \varepsilon_{\gamma\delta}$, $\varepsilon_{\gamma\delta}$ being the permutation symbol such that $\varepsilon_{12} = +1$. F^{δ} and F are the tangential and normal extrinsic forces, L^{δ} and L the corresponding couples; $S^{\gamma\delta}$ and S^{δ} are the membrane stress and cross-force resultants, $M^{\gamma\delta}$ and M^{δ} the corresponding couples. Only in special cases do (3.9) reduce to the form (2.1).

Case 6. *Equilibrium of an internally free membrane.* In all classical studies of shells, it is assumed that $L = 0$ and $M^{\delta} = 0$. If in addition we have $S^{\delta} = 0$, $L^{\gamma} = 0$, and $M^{\gamma\delta} = 0$, (3.9)₄ is satisfied identically, (3.9)₂ and (3.9)₃ reduce to

$$S^{\gamma\delta}_{,\delta} + F^{\delta} = 0, \quad S^{\gamma\delta} = S^{\delta\gamma}, \quad (3.10)$$

while (3.9)₁ becomes an algebraic equation relating $S^{\gamma\delta}$ and F :

$$b_{\gamma\delta} S^{\gamma\delta} + F = 0. \quad (3.11)$$

A tensor \mathbf{S} satisfying (3.10) and (3.11) is said to describe a *state of membrane stress* upon the surface whose fundamental forms are \mathbf{a} and \mathbf{b} . When $F^{\delta} = 0$, (3.10) is of the form (2.1). If (3.10) is solved, substitution of the result into

⁵ However, it would be permissible to use an invariant formulation throughout, since results of TOUPIN show that the laws of classical mechanics may be expressed as tensorial equations in an affine space-time such that the affine connection vanishes in all inertial frames. See § 4 of his *World invariant kinematics*, Arch. Rational Mech. Anal. **1**, 181–211 (1958). We emphasize that neither his treatment nor ours presupposes a four-dimensional metric or makes any commitment in regard to relativity.

⁶ See, e.g., § 26 of J. L. ERICKSEN & C. TRUESDELL: *Exact theory of stress and strain in rods and shells*. Arch. Rational Mech. Anal. **1**, 295–323 (1958).

(3.11) yields a differential equation to be satisfied by the stress function or potential. Solution of (3.10) by itself constitutes an intrinsic problem; when (3.11) is taken into account as well, the problem is no longer an intrinsic one.

Part I. Flat spaces

4. General solution for a flat space of n dimensions. In order that $T^{km}_{,m}=0$ in a Euclidean space of any number of dimensions, application of the classical theorem of the vector potential shows it to be necessary and sufficient that⁷

$$T^{km} = b^{km p}_{,p}, \quad \text{where} \quad b^{km p} = -b^{kp m}. \quad (4.1)$$

The condition $T^{km} = T^{mk}$ may now be written in the form

$$(b^{km p} - b^{mk p})_{,p} = 0. \quad (4.2)$$

This condition, in turn, is equivalent to the existence of a tensor c such that

$$b^{km p} - b^{mk p} = c^{km p q}_{,q}, \quad (4.3)$$

where

$$c^{km p q} = -c^{kmp q} = -c^{mk p q}. \quad (4.4)$$

Therefore

$$2b^{km p} = (c^{pkm q} + c^{mpqk} + c^{kmpq})_{,q}, \quad (4.5)$$

so that (3.1) becomes

$$T^{km} = h^{pkm q}_{,qp}, \quad (4.6)$$

where

$$\begin{aligned} h^{pkm q} &\equiv \frac{1}{2} (c^{pkm q} + c^{mqpk}), \\ h^{pkm q} &= -h^{kmp q} = -h^{pkm q} = h^{mqpk}. \end{aligned} \quad (4.7)$$

The elegant foregoing derivation, given by DORN & SCHILD⁸, shows that (4.6) furnishes the **general solution of Cauchy's laws** for equilibrium of an internally free body in a flat space of any dimension.

⁷ As observed by GWYTHER [1913], for a stress tensor which is not symmetric the analysis breaks off at this point. In a more general system of mechanics initiated by VOIGT and developed by E. & F. COSSERAT, not only is the stress tensor not generally symmetric, but also there is a couple stress tensor in addition to the usual one. In this system of mechanics, all equations of equilibrium are prescriptions of divergences. Stress functions may be introduced by the procedure given above; cf. GÜNTHER [1958, 2, § 3].

⁸ [1956, 2]. The basic idea was suggested by BELTRAMI [1892, 2]. Cf. also MORERA [1892, 3].

A variant of this derivation had been given earlier by GÜNTHER [1954, 1, § 1]. He begins by introducing the skew-symmetric dual tensor of fourth order,

$$(*) \quad T_{km p q} = e_{km n} e_{p q s} T^{ns},$$

which he interprets as a "transversal stress tensor". An explicit solution for $T_{km p q}$ in terms of stress functions is simply obtained. In DORN & SCHILD's proof, the duals appear in (5.1). While GÜNTHER's solution for the *dual* tensor is indeed valid, as he says, in a flat space of any dimension, to derive from it the ordinary stress tensor we must use the inverse of (*) and hence presume the number of dimensions to be three. Of course GÜNTHER's proof can be adjusted to the n -dimensional case also, but DORN & SCHILD's proof is equally simple and natural in all cases.

To infer (4.7), the theorem of the vector potential has been applied twice, viz., to (4.1) and to (4.2). From writings on the theory of the potential⁹ may be inferred the following **theorem of completeness**: Let \mathcal{D} be a regular region¹⁰ with boundary \mathcal{B} . In \mathcal{D} , let the components T^{km} be twice continuously differentiable functions satisfying (2.1), and in $\mathcal{D}+\mathcal{B}$, once continuously differentiable. If \mathcal{D} is infinite, assume that $T^{km}=O(r^{-2-\delta})$ and $T^{km}_{,p}=O(r^{-3-\delta})$, $\delta>0$. Then there exist infinitely many tensors h^{pkmq} such that T^{km} may be expressed in the form (4.6).

The conditions stated are merely sufficient for completeness and are by no means the weakest possible¹¹.

From (4.6) it is clear that the tensor h^{pkmq} is indeterminate to within a tensor ${}_0h^{pkmq}$ such that ${}_0h^{pkmq}_{,qp}=0$. We shall not take up the general problem of finding conditions sufficient to individualize a particular tensor within this class, though we shall allude to special cases below.

5. Three-dimensional flat space. In the three-dimensional case, set

$$a_{rs} \equiv \frac{1}{4} e_{rp} e_{sk} e_{sqm} h^{pkmq}, \quad (5.1)$$

so that

$$h^{pkmq} = e^{krp} e^{msq} a_{rs}, \quad a_{rs} = a_{sr}. \quad (5.2)$$

Then (4.6) becomes the **Gwyther-Finzi general solution**¹²:

$$T^{km} = e^{krp} e^{msq} a_{rs,pq}. \quad (5.3)$$

⁹ The result is not to be read off from any work I have seen, but a proof may be constructed by use of the apparatus assembled to prove the ordinary theorem of the vector potential. Cf., e.g., § 83 of R. S. PHILLIPS: *Vector Analysis*, New York, Wiley 1933 (the method used in § 49 yields only local completeness). More detailed consideration is given to the problem by L. LICHENSTEIN, §§ 12–14 of Chapter III of *Grundlagen der Hydromechanik*, Berlin, Springer, 1929. In these treatments, the potentials are rendered unique by requiring them to be solenoidal, but this is not convenient here.

¹⁰ Regular regions are defined in §§ 8–9 of Ch. IV of O. D. KELLOGG: *Foundations of Potential Theory*, Berlin, Springer, 1929. Such regions need not be bounded or simply connected, and their bounding surfaces may have corners and edges, subject to specified restrictions.

¹¹ I do not attempt to obtain in the present case a weakening of hypotheses so as to allow certain discontinuities in the components T^{km} and their derivatives corresponding to LICHENSTEIN's treatment of the ordinary theorem of the vector potential, because it seems preferable to base a more general analysis on the concepts of H. WEYL: *The method of orthogonal projection in potential theory*, Duke Math. J. 7, 411–444 (1940), but this requires a whole new program from the start. Furthermore, general discontinuities in T^{km} are not of interest, since the integral conservation laws to which (2.1) are equivalent for continuously differentiable fields T^{km} imply restrictions upon the discontinuities allowable. A truly general treatment should be based upon the integral conservation laws throughout.

¹² GWYTHYER [1912] obtained (5.3) in orthogonal curvilinear co-ordinates, writing out the special cases appropriate to rectangular Cartesian, cylindrical polar, and spherical polar co-ordinates (cf. also [1911]). His steps are such as to imply the necessity of the result; the sufficiency is immediate. B. FINZI [1934, 2, § 3] observed that (5.3) yields a symmetric tensor satisfying $T^{km}_{,m}=0$; for a proof of completeness he was content to refer to the previously known fact that the special cases (5.5) and (5.6) are complete. The Cartesian tensor form of (5.3) is almost evident from an equation of KLEIN & WIEGHARDT [1905, Eq. (33)], but they did not obtain it.

FINZI noticed that the tensor \mathbf{a} is indeterminate¹³ to within an arbitrary symmetric tensor ${}_0\mathbf{a}$ satisfying $e^{krs}e^{mqs}{}_0a_{rs,pq}=0$. From a known result concerning the conditions of compatibility, such a tensor is of the form ${}_0a_{mr}=b_{(m,r)}$, where \mathbf{b} is an arbitrary vector. For a given \mathbf{a} , in rectangular Cartesian co-ordinates, we may choose \mathbf{b} so as to satisfy one or the other of the conditions

$$\begin{aligned} b_{(m,r)} &= -a_{mr}, & r \neq m \\ b_{m,m} &= -a_{mm} \quad (\text{unsummed}). \end{aligned} \quad (5.4)$$

These two choices of \mathbf{b} show that for a particular rectangular Cartesian co-ordinate system there is no loss in generality in assuming in the first place that \mathbf{a} is diagonal, or that the diagonal components of \mathbf{a} are zero. The former alternative yields

$$T_{xx} = a_{,zz}^2 + a_{,yy}^3, \quad T_{xy} = -a_{,yx}^3, \quad \text{etc.}, \quad (5.5)$$

with $a^1 \equiv a_{xx}$, $a^2 \equiv a_{yy}$, $a^3 \equiv a_{zz}$; the latter alternative yields

$$T_{xx} = -2a_{,yz}^4, \quad T_{xy} = (a_{,x}^4 + a_{,y}^5 - a_{,z}^6)_{,z}, \quad \text{etc.}, \quad (5.6)$$

with $a^4 \equiv a_{23}$, $a^5 \equiv a_{31}$, $a^6 \equiv a_{12}$. These two forms of the general solution were obtained by MAXWELL and MORERA¹⁴, respectively. As follows from the special choices of \mathbf{a} made to derive them, these special forms are not invariant under transformations even of rectangular Cartesian co-ordinates. The explicit form for (5.3), with no restrictions on the six potentials, may be obtained by adding together the right-hand sides of (5.5) and (5.6). Other special choices of the potentials are possible¹⁵, but it by no means follows that a solution obtained by imposing three arbitrary conditions on the six potentials a_{km} remains complete¹⁶.

¹³ This fact is used by PERETTI [1949, 2] to show that it is possible to choose \mathbf{a} in such a way that from its components may be obtained simple expressions for the resultant force and moment of the stresses on a surface. Cf. also BLOKH [1950, 1]. GÜNTHER [1954, 1, § 3] gives simple expressions for the resultant force and torque on a body in terms of integrals of stress functions around particular curves. SCHAEFER [1955, 3] discusses the nature of ${}_0\mathbf{a}$.

¹⁴ [1868] [1869]; [1892, 1]. Using results given by BELTRAMI [1892, 2], MORERA [1892, 3] modified his derivation so as to yield (5.5) alternatively to (5.6). Cf. also GWYTHYR [1913]. A literature devoted mainly to rediscovery of known results regarding this subject has arisen recently; cf. KUZMIN [1945], WEBER [1948], MORINAGA & NÔNO [1950, 3], SCHAEFER [1953, 5], LANGHAAR & STIPPES [1954, 1], ORNSTEIN [1954, 2].

¹⁵ Cf. MORINAGA & NÔNO [1950, 3, § 3]. BLOKH [1950, 1] lists the essentially different forms that result from such special choices: 5 in rectangular Cartesian co-ordinates, 20 in general co-ordinates, 18 in cylindrical co-ordinates, 10 in cylindrical co-ordinates for rotationally symmetric problems, 19 in spherical co-ordinates.

¹⁶ What reductions are possible is not obvious. In writings on stress functions there is a deplorable custom of inferring completeness by merely counting the number of arbitrary functions. Apart from the logical gap in such inference, its danger is illustrated by the solution written down without proof of completeness by PRATELLI [1953, 4]:

$$T^{km} = e^{kpq}e^{mrs}(F_{,p}g_{qs} + H_{,p}g_{qr} + K_{,qs}g_{pr}). \quad (A)$$

While indeed a solution for any choice of the three arbitrary potentials F, H, K , it is not general. In fact, by using the expression for the product $e^{kpq}e^{mrs}$ as a determinant of δ_n^t 's, we may put (A) into the form

$$T^{km} = P_{,q}g_{,q}^{km} - P_{,k}m, \quad (B)$$

where $P \equiv F - H + K$, and it is easy to exhibit solutions of $T^{km}_{,m} = 0$ which cannot be expressed in the form (B).

While AIRY's stress function for steady plane motion subject to plane stress is formally a special case of the foregoing, a treatment which is fully two-dimensional from the start is preferable. I do not develop it here because the existing literature is adequate¹⁷.

6. Application to the axially symmetric case. If we write out (5.3) explicitly in cylindrical polar co-ordinates, at the same time supposing that all derivatives with respect to the azimuth angle are zero, for the physical components \hat{j}^k of \mathbf{T} we obtain¹⁸

$$\begin{aligned}\widehat{r\dot{r}} &= a_{,zz}^2 + \frac{1}{r} a_{,r}^3 - \frac{2}{r} a_{,z}^5, \\ \widehat{\dot{\theta}} &= a_{,rr}^3 + a_{,zz}^2 - 2a_{,rz}^5, \\ \widehat{z\dot{z}} &= a_{,rr}^2 + \frac{2}{r} a_{,r}^2 - \frac{1}{r} a_{,r}^1, \\ \widehat{r\dot{z}} &= - \left[a_{,r}^2 + \frac{1}{r} a^2 - \frac{1}{r} a^1 \right]_{,z}, \\ \widehat{r\dot{\theta}} &= \left[a_{,r}^4 - \frac{1}{r} a^4 - a_{,z}^6 \right]_{,z}, \\ \widehat{\dot{z}} &= - a_{,rr}^4 - \frac{1}{r} a_{,r}^4 + \frac{1}{r^2} a^4 + a_{,zr}^6 + \frac{2}{r} a_{,z}^6,\end{aligned}\tag{6.1}$$

where commas denote partial derivatives, and where we have set $a^1 \equiv a_{,r}$, $a^2 \equiv a_{\theta\theta}/r$, $a^3 \equiv a_{,zz}$, $a^4 \equiv a_{\theta z}/r$, $a^5 \equiv a_{,zr}$, $a^6 \equiv a_{,r\theta}/r$. The potentials a^1 , a^2 , a^3 , and a^5 occur only in the first four members of (6.1); the potentials a^4 and a^6 , only in the last two. Axially symmetric stress distributions in which $\widehat{r\dot{\theta}} = 0$, $\widehat{\dot{z}} = 0$ are often called *torsionless*; the most general stress of this kind is obtained by setting $a^4 = 0$, $a^6 = 0$ in (6.1). In any case, the six potentials may be reduced to three in a variety of ways¹⁹. For example, if we set

$$\begin{aligned}L_{,r} &\equiv a_{,r}^2 + \frac{1}{r} a^2 - \frac{1}{r} a^1, \\ M &\equiv a_{,zz}^2 + \frac{1}{r} a_{,r}^3 - \frac{2}{r} a_{,z}^5 - L_{,zz},\end{aligned}\tag{6.2}$$

¹⁷ Discovery: AIRY [1863].

Assertion of completeness: MAXWELL [1870, pp. 192–193].

First fully satisfactory treatment, including intrinsic forms and multi-valuedness in multiply connected regions: MICHELL [1900, 1].

Adjustment to steady flow subject to hydrostatic pressure: E. R. NEUMANN [1907, 1].

Interpretation in terms of resultant torque: PHILLIPS [1934, 4], SOBRERO [1935].

General curvilinear co-ordinates: B. FINZI [1934, 2, § 1].

Energetic interpretation: L. FINZI [1956, 3, § 6].

Other references: BRAHTZ [1934, 1], BATEMAN [1936] [1938], CROCCO [1950, 2].

The completeness theorem for AIRY's function does not follow as a special case of our theorem in § 4, since a plane stress field regarded from a three-dimensional viewpoint does not vanish at ∞ . Completeness is easily established directly, however.

SCHAEFER [1956, 4] approaches the plane problem as the limit of a three-dimensional one. There result, in addition to AIRY's, two stress functions representing surface loads.

¹⁸ BRDIČKA [1957, 1, § 6]. A more symmetrical expression is given by MARGUERRE [1955, 2], but his potentials are connected by a condition of compatibility, and there is no proof of completeness.

¹⁹ Cf. BLOKH [1950, 1].

then the first four members of (6.1) become²⁰

$$\begin{aligned}\widehat{r}\widehat{r} &= L_{,rr} + M, & \widehat{\vartheta}\widehat{\vartheta} &= (rM)_{,r} + L_{,zz}, \\ \widehat{z}\widehat{z} &= L_{,rr} + \frac{1}{r}L_{,r}, & \widehat{r}\widehat{z} &= -L_{,zr}.\end{aligned}\quad (6.3)$$

The function L is LOVE's *stress function for torsionless axially symmetric stress fields*. Similarly, if we put

$$W \equiv -r^2 \left[a_{,r}^4 - \frac{1}{r}a^4 - a_{,z}^6 \right], \quad (6.4)$$

the last two members of (6.1) become²¹

$$\widehat{r}\widehat{\vartheta} = -\frac{1}{r^2}W_{,z}, \quad \widehat{\vartheta}\widehat{z} = \frac{1}{r^2}W_{,r}. \quad (6.5)$$

The function W is VOIGT's *stress function for purely torsional stress fields*.

7. Application to unsteady motion subject to plane stress. We now employ the device explained under Case 4 in § 3, which rests upon the observation that the equations of plane unsteady motion subject to plane stress can be cast into the form of the equations of equilibrium of a three-dimensional body, subject to a stress T^{km} related to the plane stress $t^{\alpha\beta}$ as follows: $T^{\alpha\beta} = t^{\alpha\beta} - \rho \dot{x}^\alpha \dot{x}^\beta$, $T^{3\alpha} = T^{\alpha 3} = -\rho \dot{x}^\alpha$, $T^{33} = -\rho$, where Greek indices have the range 1, 2. This identification is possible in rectangular Cartesian inertial co-ordinate systems. The solution (5.3) is then valid for \mathbf{T} and, in these co-ordinates, is general. In writing out the result, we distinguish explicitly all time differentiations and time components. The result, derived in a special co-ordinate system, turns out to be an equation having tensorial form under general time-independent co-ordinate transformations in the plane, *viz.*

$$\begin{aligned}-\rho &= T^{33} = e^{\delta\sigma} e^{\varphi\psi} a_{\delta\varphi, \sigma\psi}, \\ -\rho \dot{x}^\gamma &= T^{\gamma 3} = e^{\gamma\delta} e^{\varphi\psi} (a'_{\delta\varphi, \psi} - a_{3\varphi, \delta\psi}), \\ t^{\gamma\lambda} - \rho \dot{x}^\gamma \dot{x}^\lambda &= T^{\gamma\lambda} = e^{\gamma\varphi} e^{\lambda\psi} (a_{33, \varphi\psi} + a''_{\varphi\psi}) - (e^{\gamma\varphi} e^{\lambda\psi} + e^{\gamma\psi} e^{\lambda\varphi}) a'_{3\psi, \varphi},\end{aligned}\quad (7.1)$$

where a prime denotes $\partial/\partial t$. In the steady case, this reduces to AIRY's solution, $-a_{33}$ being AIRY's function A . The six potentials may be reduced to three in various ways. For example, if we choose \mathbf{a} to be diagonal in a particular rectangular Cartesian co-ordinate system, we obtain²²

$$\begin{aligned}-\rho &= a_{,yy}^I + a_{,xx}^2, & \rho \dot{x} &= a_{,xt}^2, & \rho \dot{y} &= a_{,yt}^I, \\ t_{xx} - \rho \dot{x}^2 &= -A_{,yy} + a_{,tt}^2, & t_{yy} - \rho \dot{y}^2 &= -A_{,xx} + a_{,tt}^I, \\ t_{xy} - \rho \dot{x} \dot{y} &= A_{,xy}\end{aligned}\quad (7.2)$$

generalizing AIRY's solution to the case of arbitrary plane motion.

²⁰ This solution, whose completeness is easy to prove also directly from the equations of equilibrium, was obtained by LOVE [1906, § 188]. Variants are given by BRDIČKA [1957, 1, § 4].

²¹ An equivalent result in rectangular Cartesian co-ordinates was given by VOIGT [1883, pp. 297–298]. Cf. also MICHELL [1900, 2, p. 133], PHILLIPS [1934, 3].

²² A similar but not obviously identical solution is obtained by KILCHEVSKI [1953, 1].

There is literature²³ concerning the lineal case, when $t_{xy}=0$, $\dot{y}=0$, and t_{yy} does not appear in the basic equations; to obtain the results, put $A=a^I=0$.

8. Unsteady motion in three dimensions. For a Euclidean space of n dimensions, we have the solution (4.6), the completeness of which has been established. When $n=4$, letting Greek capital indices have the range 1, 2, 3, 4, we may put

$$A_{\Sigma\Lambda\Phi\Psi} \equiv \frac{1}{4} e_{\Omega\Gamma\Sigma\Lambda} e_{\Lambda\Sigma\Phi\Psi} h^{\Omega\Gamma\Lambda\Sigma}, \quad (8.1)$$

so that

$$h^{\Gamma\Lambda\Lambda\Sigma} = e^{\Lambda\Gamma\Sigma\Phi} e^{\Lambda\Sigma\Psi\Omega} A_{\Sigma\Phi\Psi\Omega}. \quad (8.2)$$

From (4.6) it follows that²⁴

$$T^{\Gamma\Lambda} = e^{\Gamma\Lambda\Sigma\Phi} e^{\Lambda\Sigma\Psi\Omega} A_{\Sigma\Phi\Psi\Omega,\Lambda\Sigma}, \quad (8.3)$$

where (4.7) implies the following conditions of symmetry for A :

$$\begin{aligned} A_{\Sigma\Phi\Psi\Omega} &= -A_{\Phi\Sigma\Psi\Omega} = -A_{\Sigma\Phi\Omega\Psi}, \\ A_{\Sigma\Phi\Psi\Omega} &= A_{\Psi\Omega\Sigma\Phi}. \end{aligned} \quad (8.4)$$

Inspection of (8.3) shows that only the second of these sets of symmetry conditions is essential; the first set may be abandoned without impairing the solution. The tensor A is indeterminate to within a tensor ${}_0A$ such that $e^{\Gamma\Lambda\Sigma\Phi} e^{\Lambda\Sigma\Psi\Omega} {}_0A_{\Sigma\Phi\Psi\Omega,\Lambda\Sigma}$. The most general such tensor is a linear combination of tensors of the type $B_{\Sigma\Phi\Psi,\Omega}$; to satisfy the essential symmetry condition (8.4)₃ we may choose $B_{\Sigma\Phi\Psi,\Omega} + B_{\Psi\Omega\Sigma,\Phi} = 0$; if, finally, we choose to satisfy also the conditions (8.4)_{1,2}, we have

$$\begin{aligned} A^0_{\Sigma\Phi\Psi\Omega} &= 4B_{[\Sigma\Phi][\Psi\Omega]} + 4B_{[\Psi\Omega][\Sigma\Phi]}, \\ &= B_{\Sigma\Phi\Psi,\Omega} - B_{\Sigma\Phi\Omega,\Psi} - B_{\Phi\Sigma\Psi,\Omega} + B_{\Phi\Sigma\Omega,\Psi} + \\ &\quad + B_{\Psi\Omega\Sigma,\Phi} - B_{\Psi\Omega\Phi,\Sigma} - B_{\Omega\Psi\Sigma,\Phi} + B_{\Omega\Psi\Phi,\Sigma}. \end{aligned} \quad (8.5)$$

To adapt this general solution for equilibrium in a flat four-dimensional space to unsteady motion in three-dimensional space, we proceed just as in § 7, except that we use (3.6)–(3.8). Writing explicitly all time components and time differentiations, we find that our result, derived in a special co-ordinate system, turns out to be in tensorial form under time-independent transformations of the space co-ordinates alone. These formulae, valid for all inertial curvilinear

²³ The first step was taken by EULER [1757, §§ 48–49]; the result was worked out by W. KIRCHHOFF [1930, 2, § 1] after being suggested by E. R. NEUMANN [1907, § 6]; it has been rediscovered by McVITTIE [1953, 3, § 3].

²⁴ B. FINZI [1934, 2, § 5] wrote down (8.3) and symmetry conditions consisting in (8.4) and a further requirement which I do not verify; he was content to infer completeness by counting the number of assignable arbitrary functions. A somewhat involved proof was given by MORINAGA & NÔNO [1950, 3, § 4]; I do not follow the argument whereby they claim [*ibid.* § 5] to establish the alternative form

$$T^{\Gamma\Lambda} = e^{\Gamma\Lambda\Phi\Psi} e^{\Lambda\Sigma\Psi\Omega} A_{\Phi\Sigma,\Lambda\Sigma};$$

they give corresponding results in n dimensions.

co-ordinate systems, are²⁵

$$\begin{aligned} -\varrho &= T^{44} = e^{smn} e^{pq\tau} A_{m n q r, s p}, \\ -\varrho \dot{x}^k &= T^{k4} = -e^{ksm} e^{pq\tau} (A'_{s m q r, p} + 2A_{4 m q r, s p}), \\ t^{km} - \varrho \dot{x}^k \dot{x}^m &= T^{km} = 4e^{kpq} e^{msn} A_{4 p 4 s, q n} + \\ &\quad + 2(e^{kpq} e^{msn} + e^{ksn} e^{mpq}) A'_{4 n p q, s} + e^{kpq} e^{msn} A''_{p q s n}, \end{aligned} \quad (8.6)$$

where a prime denotes $\partial/\partial t$, and where the symmetries of the tensors A_{pqmn} , A_{4npq} , and A_{4p4q} may be read off from (8.4). This is **Finzi's general solution** of the equations of balance of mass and momentum in an inertial frame. In the case of equilibrium, this solution reduces to (5.3) with $4A_{4p4n} = a_{pm}$.

In a particular co-ordinate system, by means of (8.5) we may impose 14 conditions upon the 20 independent potentials occurring in the solution (8.6). For example, we are tempted to take A_{4p4m} as diagonal, $A_{4mq\tau}$ as zero, and $A_{m n q r}$ as zero when $m \neq q$ and $n \neq r$. I have been unable to prove that it is possible to choose **B** in (8.5) in such a way as to justify this special choice of **A**; if it is legitimate, then, writing $A^1 \equiv 4A_{4141}, \dots, A^4 \equiv 2A_{2323}, \dots$, in this case we reduce (8.3) to the form

$$\begin{aligned} -\varrho &= A^4_{,xx} + A^5_{,yy} + A^6_{,zz}, & \varrho \dot{x} &= A^4_{,xt}, \dots \\ t_{xx} - \varrho \dot{x}^2 &= A^2_{,zz} + A^3_{,yy} + A^4_{,tt}, \dots, & t_{xy} - \varrho \dot{x} \dot{y} &= -A^3_{,xy}, \dots, \end{aligned} \quad (8.7)$$

extending MAXWELL'S formulae (8.5) so as to yield a simple yet general solution in terms of six potentials.

9. Status of the problem for flat spaces. The foregoing sections show that the *general* problem is solved in all detail for flat spaces of any dimension. It is to be hoped that this exposition may save the labor of rediscovery for those interested in the subject, enabling them to proceed at once to applications.

When theories of materials were first proposed in the eighteenth century, solutions in arbitrary functions such as those presented above were sought earnestly, but, for the most part, sought in vain. In the nineteenth century, researches on partial differential equations turned away from such general solutions so as to concentrate upon boundary-value problems. When, in the twentieth century, the general solutions were at last obtained, slight attention was paid to them, and to this day they remain virtually unknown. Though scant use has been made of them so far, they may be enlightening for studies of underdetermined systems, where the conventional viewpoint of partial differential equations has gained little.

²⁵ B. FINZI [1934, 2, § 6]. In rectangular Cartesian co-ordinates, this result is rediscovered by ARZHANIKH [1952, 1] (while he uses 21 potentials, obviously one may be eliminated). KILCHEVSKI [1953, 1] observes that if R^{Γ}_{Λ} is the Ricci tensor based on the Riemannian metric tensor $G_{\Gamma\Lambda}$, then in any Riemannian 4-space the quantities $T^{\Gamma\Lambda} \equiv R^{\Gamma\Lambda} - \frac{1}{2} G^{\Gamma\Lambda} R^{\Phi}_{\Phi}$ satisfy $T^{\Gamma\Lambda}_{;\Lambda} = 0$. Putting $G_{\Gamma\Lambda} = \delta_{\Gamma\Lambda} + \varepsilon H_{\Gamma\Lambda}$, he calculates $T^{\Gamma\Lambda} = \varepsilon Q^{\Gamma\Lambda} + O(\varepsilon^2)$; hence follows $\partial \dot{Q}^{\Gamma\Lambda} / \partial x^{\Lambda} = 0$, so that $Q^{\Gamma\Lambda}$, in rectangular Cartesian co-ordinates, gives a solution of the type presented in the text above. A similar approach, involving a detour through relativity theory, is presented by McVITTIE [1953, 3, § 2]; that his solution is not general is remarked by WHITHAM [1954, 3], who obtains what appears to be a special case of FINZI'S solution in a special co-ordinate system. Rediscoveries of other special cases are made by MILNE-THOMSON [1957, 3] and by BLANKFIELD & McVITTIE [1959, 1 & 2].

Application to the simpler theories of special materials is more difficult than might at first appear. The additional conditions imposed upon the stress functions by particular constitutive equations usually turn out to be so complicated as to baffle attempts to satisfy them. Consider, for example, the condition of plane isochoric irrotational flow subject to hydrostatic pressure. All fields of velocity and pressure of this class are included among those given by (7.1). The simple conditions $\dot{x}'_{,\nu}=0$ and $\dot{x}_{[\nu,\delta]}=0$ when applied to (7.1) yield a complicated non-linear system of equations to be satisfied by the stress functions $a_{\delta\varphi}$ and $a_{3\varphi}$; the classical formulation in terms of velocity-potential and stream-function is preferable in every way. In general, the *more* special conditions are imposed, the less practical is a solution based upon use of stress functions of the second kind.

A more promising application, initiated by FILONENKO-BORODICH²⁶ for the case of a rectangular parallelepiped, consists in determining the most general state of stress consistent with assigned values of the stress vector upon a certain boundary. Such results can lead to valuable comparisons, estimates, and mean values by which the results of different theories of materials, or of all possible theories, are compared. Cf. the remarks at the end of § 1.

Part II. Curved spaces

10. Formal intrinsic solution in an affine space. In a flat space, the similarity in form between the conditions of compatibility and the general solutions (5.4) and (5.5), or (5.3), has been remarked for half a century²⁷. This similarity may be used to yield a direct solution²⁸ of (2.1) based upon the classical principle of virtual work. Indeed, we shall need only a special case. If $T^{km}=T^{mk}$, for any vector c_k in an n -dimensional space we have

$$\int_V T^{km} c_{(k,m)} dV = \oint_S T^{km} c_k da_m - \int_V T^{km}_{,m} c_k dV. \quad (10.1)$$

²⁶ [1951, 1 & 2] [1957, 2]. The work rests on trigonometric series or special functional forms.

²⁷ Specifically, the tensor stress function is indeterminant to within a tensor satisfying the conditions of compatibility. A similar observation formed the basis from which B. FINZI [1934, 2, §§ 4, 7] conjectured his stress functions for spaces of constant curvature. See § 12, below. It follows from the results we are about to present, however, that this method of conjecture is valid only when the compatibility operator is self-adjoint, as is the case in flat spaces and in spaces of constant curvature.

²⁸ The method was given by TRUESDELL [1957, 5] for the two-dimensional case, revising work of L. FINZI [1956, 3]. Earlier SCHAEFER [1953, 5, § 4] had introduced multipliers in just the same way, concluding that „Jeder Verträglichkeitsbedingung ist eine Spannungsfunktion zugeordnet,“ but his presentation employs results of a kind valid only in flat spaces, and he did not mention any further possibilities. GÜNTHER [1954, 1, § 2] in effect noted the method and remarked that the tensor of stress functions in a three-dimensional flat space may be interpreted as Lagrangian multipliers expressing the reactions against the geometrical constraints but concluded that “no new point of view results”.

Successful treatment of the converse problem is older. As was observed by LOCATELLI [1940, 1 & 2] for a flat space and more generally by L. FINZI [1956, 3, § 10], if the general solution of the equations of equilibrium in terms of stress functions is known, the form of the conditions of compatibility may be derived from the converse of the principle of virtual work, which in an elastic context is named after MENABREA or CASTIGLIANO. It was on this basis that L. FINZI correctly conjectured for membranes the forms of solution presented in § 13, below.

This identity, being based only upon Green's theorem, is valid and invariant in any space where covariant differentiation is defined. For a fixed region V , it follows from (10.1) that a symmetric tensor T^{km} satisfies $T^{km}_{,m}=0$ in V if and only if

$$\int_V T^{km} d_{km} dv = 0 \quad (10.2)$$

for all tensors d_{km} such that the system

$$d_{km} = c_{(k,m)} \quad (10.3)$$

has a solution c_k satisfying the condition $T^{km} c_k da_m = 0$ upon S . Since we are interested in satisfying the differential equations everywhere, we may adjust the region so that the boundary condition is satisfied²⁹; in effect, we need only consider the interior condition³⁰ (10.3). In other words, the variational condition (10.2) for all d_{km} such that the system (10.3) is compatible is equivalent to the differential equation $T^{km}_{,m}=0$, provided $T^{km} = T^{mk}$.

Let it be supposed that the conditions of compatibility for (10.3), in the space considered, are of the form

$$\begin{aligned} 0 &= \mathcal{L}_{ijkl}(\mathbf{d}), \\ &= \alpha_{ijkl}{}^{pq} d_{pq} + \beta_{ijkl}{}^{pqr} d_{pq,r} + \\ &\quad + \gamma_{ijkl}{}^{pqrs} d_{pq,rs} + \delta_{ijkl}{}^{pqrst} d_{pq,rst} + \\ &\quad + \dots, \end{aligned} \quad (10.4)$$

where the sum is finite. By a standard device, this system of side conditions is taken into account by use of multipliers P^{ijkl} , so that the restricted variational condition (10.2) is replaced by the modified condition

$$\int_V [T^{km} d_{km} - P^{ijkl} \mathcal{L}_{ijkl}(\mathbf{d})] dv = 0 \quad (10.5)$$

for unrestricted \mathbf{d} . By applying Green's theorem, we put (10.5) into the form

$$\int_V [T^{km} - \bar{\mathcal{L}}^{km}(\mathbf{P})] d_{km} dv + \oint_S \dots, \quad (10.6)$$

where the dots stand for a surface integral not affecting the result we seek, and where $\bar{\mathcal{L}}$ is the operator adjoint to \mathcal{L} :

$$\bar{\mathcal{L}}^{km}(\mathbf{P}) = \alpha_{rstu}{}^{km} P^{rstu} - (\beta_{rstu}{}^{kmq} P^{rstu})_{,q} + (\gamma_{rstu}{}^{kmqv} P^{rstu})_{,qv} - \dots \quad (10.7)$$

Since (10.6) is to hold for arbitrary d_{km} , it follows that

$$T^{km} = \bar{\mathcal{L}}^{km}(\mathbf{P}). \quad (10.8)$$

²⁹ Cf. also the argument used by DORN & SCHILD [1956, 2].

³⁰ It is worth noticing that from (10.1) there follows another characterization which is not variational, as follows: For a given symmetric T^{km} such that $T^{km}_{,m}$ is continuous, let the system (10.3) be such that in any neighborhood of P there exists a region V with boundary S such that (10.3) admits n linearly independent solutions c_k which vanish upon S . Then $T^{km}_{,m}=0$ at P if and only if (10.2) holds for all such V .

The condition of symmetry $(2.1)_2$ has not been considered explicitly, but it is obviously satisfied because the tensors α, β, \dots occurring in (10.4) may be taken as symmetric in p and q without loss of generality, so that $\mathcal{L}^{km} = \mathcal{L}^{mk}$. Thus the result (10.8) always yields a symmetric \mathbf{T} . Conversely, the necessity of (10.8) follows from (10.2), and this in turn is equivalent to $(2.1)_2$ only if \mathbf{T} is assumed to be symmetric.

Thus our result (10.8) is the *general solution of the equations of motion* (2.1) in an affine n -dimensional space. The problem of solving (2.1) is thereby reduced to the problem of finding the explicit form (10.4) for the conditions of compatibility of the differential system (10.3), in the particular space considered.

When we attempt to phrase the result as a rigorous theorem subject to precise conditions of regularity, we run up against an unexpected difficulty. While the use of multipliers to eliminate constraints which are partial differential equations goes back to LAGRANGE, is standard practice in continuum mechanics, and has been used without comment by at least one eminent pure mathematician³¹, the modern literature seems to be absolutely silent regarding this subject. A straightforward attack reduces the problem to one of existence of solution of (10.8), regarded as a partial differential equation satisfied by \mathbf{P} when \mathbf{T} is given. However, there are indications that the entire approach through the principle of virtual work ought properly to be regarded in terms of a principle of invariance. Such a principle will not be sought here; for the present, we recognize that the solution (10.8) is to be regarded as formal.

11. The form of the conditions of compatibility in a Riemannian space. In the previous section, the general problem of solving the system (2.1), in an affine space, has been reduced to that of calculating explicitly the conditions of compatibility for the system (10.3) in the same space. For completeness, we indicate the nature of this problem of differential elimination and summarize the general information concerning it available at present.

Our problem is to eliminate c_k between (10.3) and its first p derivatives, for some finite p , where

$$c_{k,m} \equiv \frac{\partial c_k}{\partial x^m} - \Gamma_{km}^p c_p. \quad (11.1)$$

Regard the field c_k as given, and let d_{km} be defined by (10.3); put $w_{km} \equiv c_{[k,m]}$. For a symmetric connection $\mathbf{\Gamma}$ we have the identities

$$\begin{aligned} w_{km,p} + w_{mp,k} + w_{pk,m} &= 0, \\ c_{k,m,p} - c_{k,p,m} &= c_q R_{km}^q c_p, \\ h_{km,p,s} - h_{km,s,p} &= h_{kq} R_{mp}^q c_s + h_{qm} R_{kp}^q c_s. \end{aligned} \quad (11.2)$$

Since $c_{k,mq} = d_{km,p} + w_{km,p}$, by forming $c_{k,m,p} - c_{k,p,m}$ and using (11.2)_{1,2} we find that

$$w_{mp,k} = d_{km,p} - d_{kp,m} - R_{km}^q c_q. \quad (11.3)$$

³¹ E. HELINGER, §§ 3d and 7e of: Die allgemeinen Ansätze der Mechanik der Kontinua. Enz. Math. Wiss. IV 2.2, H. 5, 601–694 (1913): „so kann man nach dem bekannten Kalkül der Variationsrechnung ... zugehörige Lagrangesche Faktoren ... einführen ...“.

Now forming $w_{m p, k s} - w_{m p, s k}$ and using (11.2)₃, we obtain the identity³²

$$\begin{aligned} d_{k m, p s} - d_{k p, m s} - d_{s m, p k} + d_{s p, m k} + \\ + (R_{s m p, k}^q - R_{k m p, s}^q) c_q + R_{s m p}^q d_{q k} - R_{k m p}^q d_{q s} - \\ - R_{k m p}^q w_{q s} + R_{s m p}^q w_{q k} + R_{p k s}^q w_{q m} - R_{m k s}^q w_{q p} = 0. \end{aligned} \quad (11.4)$$

For a flat space, (11.4) reduces to the Kirchhoff-St. Venant conditions

$$d_{k m, p s} - d_{k p, m s} - d_{s m, p k} + d_{s p, m k} = 0. \quad (11.5)$$

In more general spaces, we are to derive from (11.4) and from the geometry of the space further identities, by means of which c_q and $w_{k m}$ are to be eliminated.

PALATINI³³ has indicated how the calculation can be initiated in the case of a metric space. First, we replace sums of the type $R^q \dots a_q$ by corresponding sums $R_q \dots a^q$, then use the relation $R_{k m p s} = R_{p s k m}$ and the Bianchi identities so as to obtain $(R_{s m p, k}^q - R_{k m p, s}^q) c_q = -R_{m p s k, q} c^q$. Thus (11.4) becomes

$$\begin{aligned} \delta_{k p}^{r q} \delta_{m s}^{t u} d_{r t, q u} - R_{m p s k, q} c^q - R_{m p s q} d_{k}^q + R_{m p k q} d_{s}^q + \\ + R_{m p k q} w_{q s} - R_{m p s q} w_{q k} - R_{k s p q} w_{q m} + R_{k s m q} w_{q p} = 0. \end{aligned} \quad (11.6)$$

Raising the index k and then contracting upon k and m yields³⁴

$$I_{, p s} + d_{p s, k}^k - 2d_{(p, s) q}^q - 2d_{(p}^q R_{s) q} - 2w_{(p}^q R_{s) q} - R_{p s, q} c^q = 0, \quad (11.7)$$

where (11.2)₃ has been used, where $R_{p q} \equiv R_{p q h}^h$, and where $I \equiv d_k^k$. Raising the index p in (84.10) and then contracting upon q and s yields³⁵

$$I_{, p}^p d_{, k p}^k - R_{p q} d^{p q} = R_{, q} c^q \quad (11.8)$$

where $R \equiv R_k^k$.

Equations (11.6) are to be regarded as a system of $\frac{1}{2}n(n+1)$ non-homogeneous linear equations for the $\frac{1}{2}n(n+1)$ unknowns c^q and $w^{k m}$ in terms of the tensor \mathbf{d} and its second derivatives. Substituting the solution of these equations into (11.6) yields, in principle, a set of differential equations for \mathbf{d} alone³⁶. The geometrical meaning of the condition of solubility expressed in terms of a determinant formed from the components $R_{p q}$ and $R_{p q, s}$ is not known; of course, it is not satisfied in a flat space, nor is it relevant, since the appropriate condition has already been shown to be (11.5).

More generally, directly from (10.3) we see that the nature of the solution depends very strongly upon the geometry of the space considered. In spaces

³² RICCI [1888, Eq. (20)].

³³ [1934, 3]. The method was sufficiently indicated by his earlier analysis of a space of constant curvature [1916]. Cf. also ANDRUETTO [1932, 1 and 2], AGOSTINELLI [1933], GRAIFF [1958, 1, §§ 6–9].

³⁴ This identity was given in a special co-ordinate system by PALATINI [1934, 3]; an equivalent general form is due to GRAIFF [1958, 1, Eq. (22)].

³⁵ GRAIFF [1958, 1, Eq. (23)].

³⁶ This seems to be a correct substitute for the method of PALATINI [1934, 3], who notes that it is possible to choose co-ordinates so that at a fixed point (11.7) with $p = s$ becomes a system of n linear equations for the n components c^q but fails to note that the solution of this system is not sufficient, in general, to obtain the derivatives of c^q and so to determine $w^{k m}$.

admitting a group of "motions", *i.e.*, instantaneously rigid velocity fields, the solution of (10.3), if compatible, is determinate only to within such a motion. In spaces where instantaneously rigid motions are not possible, the system (10.3) will determine \mathbf{c} uniquely from a given \mathbf{d} , if compatible. Signorina GRAIFF³⁷ has initiated a discussion of the problem in terms of the principal invariants of the tensor R_{pq} , according as these are linearly independent and non-constant or not. The details are complex, and the problem must still be regarded as open.

In a space of constant curvature, we have

$$\begin{aligned} n(n-1) R_{km}^q &= R(g_{km} \delta_p^q - g_{kp} \delta_m^q) \\ n R_{km} &= R_{qkm}, \end{aligned} \quad (11.9)$$

where $R \equiv R_k^k = \text{const.}$ The terms involving \mathbf{c} and \mathbf{w} in (11.6) vanish. The resulting conditions, easily seen to be sufficient as well as necessary, are³⁸

$$\delta_{kp}^{rq} \delta_{ms}^{tu} (d_{r,t,qu} - \frac{R}{n(n-1)} d_{r,t} g_{qu}) = 0. \quad (11.10)$$

When $n=2$, there is but one linearly independent non-vanishing component of (11.10), which may be written in the equivalent forms

$$\begin{aligned} e^{\alpha\beta} e^{\gamma\delta} \bar{d}_{\alpha\gamma,\beta\delta} + K I &= 0, \\ I_{,\alpha}^\alpha - \bar{d}^{\alpha\beta}_{,\alpha\beta} + K I &= 0, \end{aligned} \quad (11.11)$$

where K is the total curvature, $K = -\frac{1}{2}R$.

More generally, when $n=2$ and K is arbitrary we have $R_{\alpha\beta\gamma\delta} = K e_{\alpha\beta} e_{\gamma\delta}$ and $R_{\alpha\beta} = -K a_{\alpha\beta}$, where \mathbf{a} is the metric tensor. Again the terms containing \mathbf{w} in (11.6) are annulled, and we have

$$e^{\alpha\beta} e^{\gamma\delta} d_{\alpha\gamma,\beta\delta} + K I = -2K_{,\alpha} c^\alpha. \quad (11.12)$$

Further differentiations are needed to eliminate \mathbf{c} . B. FINZI³⁹ has shown that for a surface applicable upon a surface of revolution the final condition assumes the form

$$(K_{,\gamma} K^{,\gamma} a^{\alpha\beta} + K^{,\alpha} K^{,\beta}) d_{\alpha\beta} + K K^{,\gamma} \bar{d}^{\alpha\beta}_{,\alpha\beta,\gamma} + e^{\alpha\beta} e^{\gamma\delta} K^{,\epsilon} d_{\beta\delta,\alpha\gamma\epsilon} = 0. \quad (11.13)$$

For a general surface, he obtained a system of three conditions of compatibility, each being a differential equation of fourth order. I remarked⁴⁰ that they must be equivalent to a single condition of fifth order. By a method of infinitesimal variation applied to a complete set of differential invariants of the surface in question, Signorina GRAIFF⁴¹ has given a definitive treatment of the problem. She has found two identities relating the quantities occurring in B. FINZI's conditions, but she has not effected the elimination explicitly. She has charac-

³⁷ [1958, 1, § 10].

³⁸ PALATINI [1916], [1934, 3], ANDRUETTO [1932, 2], FINZI [1934, 2] including a simplification when $n=3$.

³⁹ [1930, 1].

⁴⁰ [1957, 6]; the reasoning is presented in § 13, below.

⁴¹ [1957, 3].

terized the case when the least possible order of the single scalar condition is 4; in this case, as for a surface of revolution, the lines $K = \text{const.}$ are geodesic parallels, but $K_{,\alpha}^{\alpha}$ is not a function of K only.

Even in the cases when $n=3$ and $n=4$, the possibilities are various⁴², and the details are too intricate to summarize here. The fullest information now available is included in the two papers of Signorina GRAIFF already cited.

12. Application to spaces of constant curvature. In a space of constant curvature, from (11.10) it follows that the coefficients in the form (10.4) have the special forms

$$\begin{aligned}\alpha_{ijkl}{}^{pq} &= -\frac{R}{n(n-1)} \delta_{ij}^p \delta_{kl}^q g_{uv}, \\ \gamma_{ijkl}{}^{pqrs} &= \delta_{ij}^p \delta_{kl}^q \delta_{rs}^s,\end{aligned}\quad (12.1)$$

while all other coefficients vanish. The operator \mathcal{L} is self-adjoint, and from (10.8) we have the following *general solution for a space of constant curvature*:

$$T^{km} = h^{pkmq}{}_{,pq} - \frac{R}{n(n-1)} h^{qkm}{}_q, \quad (12.2)$$

where we have set

$$h^{pkmq} \equiv \delta_{rs}^p \delta_{tu}^q P^{rstu}, \quad (12.3)$$

P being the tensor in terms of which the general solution of § 10 is expressed. From (12.3), it is immediately plain that h satisfies the conditions of symmetry (4.7)_{2,3,4}. Thus the solution (4.6) for flat spaces is included in (12.2) as the special case when $R=0$.

When $n=3$, we may introduce the tensor a dual to h according to (5.1). The solution (12.2) now becomes⁴³

$$\begin{aligned}T^{km} &= e^{krs} e^{msq} (a_{rs,pq} - \tfrac{1}{6} R a_{rs} g_{pq}), \\ &= e^{krs} e^{msq} a_{rs,pq} - \tfrac{1}{6} R (g^{km} a_q^q - a^{km}),\end{aligned}\quad (12.4)$$

generalizing (5.3).

When $n=4$, we may introduce the tensor A dual to h according to (8.1) and so obtain from (12.2) the solution⁴⁴

$$T^{rA} = e^{rA\Sigma\Phi} e^{A\varepsilon\Psi\Omega} (A_{\Sigma\Phi\Psi\Omega,\varepsilon A} - \tfrac{1}{12} R A_{\Sigma\Phi\Psi\Omega} g_{\varepsilon A}), \quad (12.5)$$

generalizing (8.3). However, this form of the solution does not seem to offer any advantage over (12.2).

13. Application to the classical theory of membranes. We now consider the membrane problem: To find the general solution of (3.10), where covariant differentiation is based upon a given two-dimensional Riemannian metric a ,

⁴² GRAIFF [1958, 1, §§ 13–14].

⁴³ B. FINZI [1934, 2, § 4].

⁴⁴ B. FINZI [1934, 2, § 7] gives the solution

$$T^{rA} = e^{rA\Sigma\Phi} e^{A\varepsilon\Psi\Omega} A_{\Sigma\Phi\Psi\Omega,A\varepsilon} + K(A_{\Sigma\Phi}{}^{\Sigma\Phi} g^{rA} - 2A^r{}_{\Phi}{}^A{}_{\Phi});$$

this may be equivalent to (12.5) but is not obviously so.

and where we take $F^\delta = 0$. By direct and laborious reduction, STORCHI⁴⁵ has obtained a general solution in geodesic co-ordinates. This solution involves derivatives of orders up to 5 of a single stress function. At the present writing, a corresponding invariant solution is not yet known. In order to find it by our method, we should have to perform explicitly the reduction of B. FINZI's conditions of compatibility to a single scalar equation.

We rest content with recording the simple cases that follow effortlessly from the results already given. When the condition of compatibility is a single scalar condition $\mathcal{L}(\mathbf{d}) = 0$, the analysis proceeds just as in § 10, except that the four indices $ijkl$ are suppressed, and the multiplier reduces to a single scalar A .

For a surface of constant curvature, by comparing (11.10) with the reduced form of (10.4) we see that

$$\alpha^{\sigma\varphi} = K a^{\sigma\varphi}, \quad \beta^{\sigma\varphi\varrho} = 0, \quad \gamma^{\sigma\varphi\varrho\tau} = e^{\sigma\varrho} e^{\varphi\tau}, \quad (13.1)$$

and all other coefficients vanish. Substitution into the reduced form of (10.8) yields the *general solution for surfaces of constant curvature*⁴⁶:

$$T^{\alpha\beta} = e^{\alpha\gamma} e^{\beta\varphi} A_{,\gamma\varphi} + K A a^{\alpha\beta}. \quad (13.2)$$

Of course, this same result follows from (12.2) and the observation that when $n=2$, the most general tensor h having the symmetries (4.7)_{2,3,4} is of the form $h^{\gamma\alpha\beta\delta} = e^{\alpha\gamma} e^{\beta\delta} A$ for some scalar A .

For a surface applicable upon a surface of revolution, by comparing (11.12) with the reduced form of (10.4) we see that

$$\begin{aligned} \alpha^{e\varphi} &= K_{,\sigma} K^{,\sigma} a^{e\varphi} + K^{,e} K^{,\varphi}, \\ \beta^{\varphi\varrho\sigma} &= K K^{,\sigma} a^{e\varphi}, \quad \gamma^{e\varphi\sigma\tau} = 0, \\ \delta^{e\varphi\sigma\tau\nu} &= -e^{\tau(e} e^{\varphi)\sigma} K^{,\nu}, \end{aligned} \quad (13.3)$$

and all other coefficients vanish. Substitution into the reduced form of (10.8) yields the *general solution for a surface applicable upon a surface of revolution*

⁴⁵ [1950, 5]. STORCHI [1950, 4] observed also that if for an arbitrary surface we choose co-ordinates x, y so that $ds^2 = \lambda(dx^2 + dy^2)$, then a solution of (2.10) when $\mathbf{F} = 0$ is furnished by the formulae $\lambda^2 S^{xx} = \partial^2 A / \partial y^2$, $\lambda^2 S^{xy} = -\partial^2 A / \partial x \partial y$, $\lambda^2 S^{yy} = \partial^2 A / \partial x^2$, provided that $\partial^2 A / \partial x^2 + \partial^2 A / \partial y^2 = 0$. For such solutions the mean pressure vanishes: $\lambda(S^{xx} + S^{yy}) = S^\delta_\delta = 0$, but it is not shown that all solutions such that $S^\delta_\delta = 0$ are included. STORCHI treated the case of a surface of revolution in [1949, 3]; in [1952, 2] and [1953, 6], the case when a general solution involving derivatives of orders no higher than the fourth is possible. Cf. the corresponding condition of compatibility mentioned at the end of § 11. A solution for minimal surfaces is initiated by COLONNETTI [1956, 1]. I do not attempt to present the solutions, some classical and some recent, for special cases defined by conditions of inextensibility. For Riemannian spaces of 2, 3, and 4 dimensions, a special solution containing an arbitrary function of the total curvature is obtained by STORCHI [1957, 4].

⁴⁶ B. FINZI [1934, 2, § 2] verified that (13.2) satisfies (2.1) when K is constant but did not prove the completeness of this solution; his result is rediscovered by LANGHAAR [1953, 2].

of non-constant curvature⁴⁷:

$$\begin{aligned} T^{\alpha\beta} = & [-K K_{,\gamma}^{\gamma} a^{\alpha\beta} + K_{,\alpha}^{\alpha} K_{,\beta}^{\beta} + e^{\delta(\alpha} \epsilon^{\beta)\gamma} K_{,\gamma}^{\gamma}{}_{,\delta\epsilon}] A + \\ & + [-K K_{,\gamma}^{\gamma} a^{\alpha\beta} + e^{\epsilon(\alpha} e^{\beta)\delta} K_{,\delta}^{\gamma}{}_{,\epsilon} + 2e^{\gamma(\alpha} e^{\beta)\delta} K_{,\delta}^{\epsilon}{}_{,\gamma\epsilon}] A_{,\gamma} + \\ & + [2e^{\gamma(\alpha} e^{\beta)\delta} K_{,\delta}^{\epsilon}{}_{,\gamma\epsilon} + e^{\epsilon(\alpha} e^{\beta)\gamma} K_{,\lambda}^{\lambda}{}_{,\gamma\epsilon}] A_{,\gamma\epsilon} + \\ & + e^{\delta(\alpha} \epsilon^{\beta)\gamma} K_{,\gamma}^{\epsilon}{}_{,\delta\epsilon} A_{,\gamma\delta\epsilon}. \end{aligned} \quad (13.4)$$

[Note added in proof, September 15, 1959. The elegant paper of GRAIFF [1959, 3] on this problem has reached me too late for a full description to be included here. Her essential idea is to begin with an invariant decomposition of the three-dimensional space of symmetric covariant tensors defined upon the surface. The general case is defined as that in which the scalars K and $H \equiv K_{,\alpha} K^{\alpha}$ are functionally independent. Then the three symmetric covariant tensors $K_{,\alpha} K_{,\beta}$, $H_{,\alpha} H_{,\beta}$, and $H_{(\alpha} K_{\beta)}$ form a basis for the symmetric covariant tensors $F_{\alpha\beta}$:

$$F_{\alpha\beta} = A K_{,\alpha} K_{,\beta} + B H_{(\alpha} K_{\beta)} + C H_{,\alpha} H_{,\beta}, \quad (13.5)$$

where the scalar functions A, B, C are uniquely determined joint invariants of F, K , and H which are easy to exhibit. After calculation of $F_{\alpha\beta}{}^{,\beta}$ it turns out that the most general symmetric tensor $F_{\alpha\beta}$ such that

$$F_{\alpha\beta}{}^{,\beta} \propto K_{,\alpha} \quad (13.6)$$

corresponds to a special choice of A in (13.5). Now let φ be any scalar function, and put

$$\Phi_{\alpha\beta} = a_{\alpha\beta}(\varphi_{,\gamma}{}^{\gamma} + K \varphi) - \varphi_{,\alpha\beta}; \quad (13.7)$$

then

$$\Phi_{\alpha\beta}{}^{,\beta} = \varphi K_{,\alpha}. \quad (13.8)$$

Hence the most general solution of (2.1) is of the form $T_{\alpha\beta} = D_{\alpha\beta} + \Phi_{\alpha\beta}$ where $D_{\alpha\beta}$ is the most general solution of

$$D_{\alpha\beta}{}^{,\beta} = -\varphi K_{,\alpha} \quad (13.9)$$

for any given φ . Matching (13.6) and (13.9) determines a unique value for φ . The result is an explicit formula for the general solution of (2.1) in terms of derivatives up to fourth order of two scalar potentials. This is a result of the type to be expected in view of the conditions of compatibility discussed in § 11.

However, GRAIFF infers that *either one of the two scalars may always be taken as zero*. This contradicts the results of STORCHI, cited in footnote 45, according to which it is only an exceptional class of surfaces for which fourth derivatives of a single scalar suffice. In GRAIFF's work, a solution χ to her equation (23)₁ with $B' = 0$ or (23)₂ with $C' = 0$, for any B and C , respectively, is assumed, but the existence of solution to this problem on a curved surface is not proved.]

An essentially different approach to the membrane problem has been initiated by PUCHER⁴⁸. He combines (3.10) and (3.11) from the start, so that the intrinsic problem is never solved, and he uses special oblique co-ordinates on the surface. While a simpler derivation of his results has been given by L. FINZI⁴⁹, it does not seem easy to put the analysis into invariant form.

14. Status of the problem for curved spaces. While the powerful method of § 10 reduces the problem to one of differential elimination, the necessary

⁴⁷ Given in geodesic co-ordinates, with a proof of completeness, by STORCHI [1949, 3]; in the above form, with an imperfect proof of completeness, by L. FINZI [1956, 3].

⁴⁸ [1934, 5] [1939].

⁴⁹ [1955, 1].

calculations seem to be difficult. For physical interpretation, curved spaces are most interesting in two special cases: ordinary surfaces and four-dimensional space-time.

The most useful curved surfaces are included in the results given in § 13; the one further elimination needed for the general surface is known to be possible and doubtless will soon be effected.

No special simplification is apparent, however, when $n=4$. Furthermore, recent studies of the foundations of mechanics show that the most interesting result would be that for a not necessarily Riemannian affine space. While RICCI's identity (11.4) is valid for such a space, none of the further calculations in § 11 are applicable. Since the problem as stated is an affine one, it seems that the entry so far used is not really cogent, and that a more penetrating study of the structure of affine spaces is needed.

Acknowledgment. I am grateful to Professor ERICKSEN for assistance in preparing § 8 and to Professor STERNBERG for criticism of an earlier version.

This report was written for the Rheology Section, National Bureau of Standards, Washington, D. C.

References

- 1757 EULER, L.: Continuation des recherches sur la théorie du mouvement des fluides. *Mém. acad. sci. Berlin* [11], 316–361 (1755) = *Opera omnia* (2) 12, 92–132.
- 1863 AIRY, G. B.: On the strains in the interior of beams. *Phil. trans. r. soc. London* 153, 49–80. Abstract in *Rep. Brit. Assn.* 1862, 82–86 (1863).
- 1868 MAXWELL, J. C.: On reciprocal diagrams in space, and their relation to AIRY's function of stress. *Proc. London Math. Soc.* (1) 2, 58–60 (1865–1869) = *Papers* 2, 102–104.
- 1870 MAXWELL, J. C.: On reciprocal figures, frames, and diagrams of forces. *Trans. r. soc. Edinburgh* 26, 1–40 (1869–1872) = *Papers* 2, 161–207.
- 1883 VOIGT, W.: Allgemeine Formeln für die Bestimmung der Elasticitätsconstanten von Krystallen durch die Beobachtung der Biegung und Drillung von Prismen. *Annalen der Phys.* (2) 16, 273–321, 398–416.
- 1888 RICCI, G.: Delle derivazioni covarianti e controvarianti e del loro uso nella analisi applicata. *Studi univ. Padova comm. ottavo centenar. Univ. Bologna* 3, No. 12, 23 pp.
- 1892 1. MORERA, G.: Soluzione generale delle equazioni indefinite dell'equilibrio di un corpo continuo. *Rend. Lincei* (5) 1, 137–141.
2. BELTRAMI, E.: Osservazioni sulla nota precedente. *Rend. Lincei* (5) 1, 141–142 = *Opere* 4, 510–512.
3. MORERA, G.: Appendice alla Nota: sulla soluzione più generale delle equazioni indefinite dell'equilibrio di un corpi continuo. *Rend. Lincei* (5) 1, 233–234.
- 1900 1. MICHELL, J. H.: On the determination of stress in an elastic solid, with applications to the theory of plates. *Proc. London Math. Soc.* 31, 100–124 (1899).
2. MICHELL, J. H.: The uniform torsion and flexure of incomplete tores, with application to helical springs. *Proc. London Math. Soc.* 31, 130–146 (1899).
- 1905 KLEIN, F., & K. WIEGHARDT: Über Spannungsflächen und reziproke Diagramme, mit besonderer Berücksichtigung der Maxwellschen Arbeiten. *Arch. Math. Phys.* (3) 8, 1–10, 95–119.
- 1906 LOVE, A. E. H.: *A Treatise on the Mathematical Theory of Elasticity*, 2nd ed. Cambridge.
- 1907 NEUMANN, E. R.: Über eine neue Reduktion bei hydrodynamischen Problemen. *J. reine angew. Math.* 132, 189–215.

- 1911 GWYTHYR, R. R.: The conditions that the stresses in a heavy elastic body should be purely elastic stresses. *Mem. Manchester lit. phil. soc.* **55**, No. 20 (12 pp.).
- 1912 GWYTHYR, R. R.: The formal specification of the elements of stress in Cartesian, and in cylindrical and spherical polar coordinates. *Mem. Manchester lit. phil. soc.* **56**, No. 10 (13 pp.).
- 1913 GWYTHYR, R. F.: The specification of the elements of stress, Part II. A simplification of the specification given in Part I. *Mem. Manchester lit. phil. soc.* **57**, No. 5 (4 pp.).
- 1916 PALATINI, A.: Sulle quadriche di deformazione per gli spazi S_3 . *Atti ist. Veneto* **76**, 125–148.
- 1930 1. FINZI, B.: Sopra il tensore di deformazione di un velo. *Ist. Lombardo rend.* (2) **63**, 975–982.
2. KIRCHHOFF, W.: Reduktion simultaner partiellen Differentialgleichungen bei hydrodynamischen Problemen. *J. reine angew. Math.* **164**, 183–195.
- 1932 1. ANDRUETTO, G.: Le formule di SAINT-VENANT per gli spazi curvi a tre dimensioni. *Rend. accad. Lincei* (6) **15**, 214–218.
2. ANDRUETTO, G.: Le formule di SAINT-VENANT per le varietà V_n a curvatura costante. *Rend. accad. Lincei* (6) **15**, 792–797.
- 1933 AGOSTINELLI, C.: Le condizioni di SAINT-VENANT per le deformazioni di una varietà riemanniana generica. *Rend. accad. Lincei* (6) **18**, 529–533; (6) **19**, 22–26 (1934).
- 1934 1. BRAHTZ, J.: Notes on the AIRY stress function. *Bull. Amer. math. soc.* **40**, 427–430.
2. FINZI, B.: Integrazione delle equazioni indefinite della meccanica dei sistemi continui. *Rend. accad. Lincei* (6) **19**, 578–584, 620–623.
3. PALATINI, A.: Sulle condizioni di SAINT-VENANT in una V_n qualsivoglia. *Rend. accad. Lincei* (6) **19**, 466–469.
4. PHILLIPS, H. B.: Stress functions. *J. math. phys. M. I. T.* **13**, 421–425.
5. PUCHER, A.: Über den Spannungszustand in gekrümmten Flächen. *Beton u. Eisen* **33**, 298–304.
- 1935 SOBRERO, L.: Del significato meccanico della funzione de AIRY. *Ricerche di ingegn.* **3**, 77–80 = *Rend. Lincei* (6) **21**, 264–269.
- 1936 BATEMAN, H.: Progressive waves of finite amplitude and some steady motions of an elastic fluid. *Proc. nat. acad. sci. U.S.A.* **22**, 607–619.
- 1938 BATEMAN, H.: The lift and drag functions for an elastic fluid in two dimensional irrotational flow. *Proc. nat. acad. sci. U.S.A.* **24**, 246–251.
- 1939 PUCHER, A.: Über die Spannungsfunktion beliebig gekrümmter dünner Schalen. *Proc. 5th. int. congr. appl. mech.*, Cambridge 1938, 134–139.
- 1940 1. LOCATELLI, P.: Sulla congruenza delle deformazioni. *Rend. ist. Lombardo* **13** = (3) **4**, 451–464 (1939–1940).
2. LOCATELLI, P.: Sul principio di MENABREA. *Boll. un. mat. Ital.* (2) **2**, 342–347.
- 1945 KUZMIN, R. O.: On MAXWELL's and MORERA's formulae in the theory of elasticity. *C. r. (Doklady) acad. sci. URSS* **49**, 326–328.
- 1948 WEBER, C.: Spannungsfunktionen des dreidimensionalen Kontinuums. *Z. angew. Math. Mech.* **28**, 193–197.
- 1949 1. FINZI, B., & M. PASTORI: Calcolo tensoriale e applicazioni. Bologna, Zanichelli.
2. PERETTI, G.: Significato del tensore arbitrario che interviene nell' integrale generale delle equazioni della statica dei continui. *Atti sem. mat. fis. univ. Modena* **3**, 77–82.
3. STORCHI, E.: Integrazione delle equazioni indefinite della statica dei sistemi continui su una superficie di rotazione. *Rend. accad. Lincei* (7) **7**, 227–231.
- 1950 1. BLOKH, V.: Stress functions in the theory of elasticity [in Russian]. *Prikl. Mat. Mekh.* **14**, 415–422.
2. CROCCO, L.: On a kind of stress-function for the study of non-isentropic two-dimensional motion of gases. *Proc. 7th int. congr. appl. mech. London 1948*, **2**, 315–329.

3. MORINAGA, K., & T. NÔNO: On stress-functions in general coordinates. *J. sci. Hiroshima Univ. A* **14**, 181—194.
4. STORCHI, E.: Sulle equazioni indefinite della statica delle membrane tese su generiche superficie. *Rend. accad. Lincei* (7) **8**, 116—120.
5. STORCHI, E.: Integrazione delle equazioni indefinite della statica dei veli tesi su una generica superficie. *Rend. accad. Lincei* (7) **8**, 326—331.
- 1951 1. FILONENKO-BORODICH, M. M.: The problem of the equilibrium of an elastic parallelepiped subject to assigned loads on its boundaries [in Russian]. *Prikl. Mat. Mekh.* **15**, 137—148.
2. FILONENKO-BORODICH, M. M.: Two problems on the equilibrium of an elastic parallelepiped [in Russian]. *Prikl. Mat. Mekh.* **15**, 563—574.
- 1952 1. ARZHANIKH, I. C.: Tensor functions of hydrodynamical stresses [in Russian]. *Doklad. Akad. Nauk SSSR* **83**, 195—198.
2. STORCHI, E.: Le superficie eccezionali nella statica delle membrane. *Revista mat. univ. Parma* **3**, 339—360.
- 1953 1. KILCHEVSKI, N. A.: Stress, velocity, and density functions in static and dynamic problems in the mechanics of continuous media [in Russian]. *Doklad. Akad. Nauk SSSR* **92**, 895—898.
2. LANGHAAR, H.: An invariant membrane stress function for shells. *J. appl. mech.* **20**, 178—182.
3. McVITTIE, G. C.: A method of solution of the equations of classical gas dynamics using EINSTEIN's equations. *Q. appl. math.* **11**, 327—336.
4. PRATELLI, A.: Sulla stazionarietà di significitavi integrali nella meccanica dei continui. *Rend. ist. Lombardo* (3) **17** (86), 714—724.
5. SCHAEFER, H.: Die Spannungsfunktionen des dreidimensionalen Kontinuums und des elastischen Körpers. *Z. angew. Math. Mech.* **33**, 356—362.
6. STORCHI, E.: Sulle membrane aventi comportamento meccanico eccezionale. *Ist. Lombardo rend.* (3) **17** (86), 462—483.
- 1954 1. GÜNTHER, W.: Spannungsfunktion und Verträglichkeitsbedingungen der Kontinuumsmechanik. *Abh. Braunschweig. Wiss. Ges.* **7**, 107—112.
2. LANGHAR, H., & M. STIPPES: Three-dimensional stress functions. *J. Franklin inst.* **258**, 371—382.
3. ORNSTEIN, W.: Stress functions of MAXWELL and MORERA. *Q. appl. math.* **2**, 198—201.
4. WHITHAM, G. B.: A note on a paper by G. C. McVITTIE. *Q. appl. math.* **12**, 316—318.
- 1955 1. FINZI, L.: Sulle equazioni di PUCHER nell'equilibrio delle strutture a guscio. *Rend. ist. Lombardo* (3) **88**, 907—916.
2. MARGUERRE, K.: Ansätze zur Lösung der Grundgleichungen der Elastizitätstheorie. *Z. angew. Math. Mech.* **35**, 242—263.
3. SCHAEFER, H.: Die Spannungsfunktion einer Dyname. *Abh. Braunschweig. Wiss. Ges.* **7**, 107—112.
4. SCHAEFER, H.: Die drei Spannungsfunktionen des zweidimensionalen ebenen Kontinuums. *Österr. Ing.-Arch.* **10**, 267—277. [I am unable to see the connection found by PRAGER, *Math. Rev.* **18**, 613 (1957), between this work and one by FOX & SOUTHWELL.]
- 1956 1. COLONNETTI, G.: L'équilibre des voiles minces hyperstatiques (Le cas des voiles de surface minimum). *C. r. acad. sci. Paris* **243**, 1087—1089, 1701—1704.
2. DORN, W. S., & A. SCHILD: A converse to the virtual work theorem for deformable solids. *Q. appl. math.* **14**, 209—213.
3. FINZI, L.: Legame fra equilibrio e congruenza e suo significato fisico. *Rend. accad. Lincei* (8) **20**, 205—211, 338—342.
- 1957 1. BRDIČKA, M.: On the general form of the BELTRAMI equation and PAPKOVICH's solution of the axially symmetric problem of the classical theory of elasticity. *Czechosl. J. phys.* **7**, 262—274.
2. FILONENKO-BORODICH, M. M.: On the problem of LAMÉ for the parallelepiped in the general case of surface loads [in Russian]. *Prikl. Mat. Mekh.* **21**, 550—559.

3. GRAIFF, F.: Sulle condizioni di congruenza per una membrana. Rend. ist. Lombardo (A) **92**, 33—42.
4. MILNE-THOMSON, L. M.: A general solution of the equations of hydrodynamics. J. Fluid. Mech. **2**, 88.
5. STORCHI, E.: Una soluzione delle equazioni indefinite della meccanica dei continui negli spazi riemanniani. Rend. ist. Lombardo **90**, 369—378 (1956).
6. TRUESDELL, C.: General solution for the stresses in a curved membrane. Proc. nat. acad. sci. U.S.A. **43**, 1070—1072.
- 1958 1. GRAIFF, F.: Sulle condizioni di congruenza per deformazioni anche finite. Rend. accad. Lincei (8) **24**, 415—422.
2. GÜNTHER, W.: Zur Statik und Kinematik des Cosseratschen Kontinuums. Abh. Braunsch. Wiss. Ges. **10**, 195—213.
- 1959 1. BLANKFIELD, J., & G. C. McVITTIE: EINSTEIN'S equation and classical hydrodynamics. Arch. Rational Mech. Anal. **2**, 337—354.
2. BLANKFIELD, J., & G. C. McVITTIE: A method of solution of the equations of magneto-hydrodynamics. Arch. Rational Mech. Anal. **2**, 411—422.
3. GRAIFF, F.: Soluzione tensoriale generale delle equazioni indefinite di equilibrio per una membrana. Rend. accad. Lincei (8) **26**, 189—196.

Appendix

Bibliography of works on stress functions for linear elasticity and related theories, not repeating items already cited as references to the foregoing article.

- 1829 E 1. POISSON, S. D.: Addition au mémoire sur l'équilibre et le mouvement des corps élastiques. Mém. Acad. Sci. Paris **8**, 623—627.
- 1851 E 1. STOKES, G. G.: On the dynamical theory of diffraction. Trans. Cambridge Phil. Soc. **9** (1851—1856), 1—62 = Papers **2**, 243—328.
- 1852 E 1. LAMÉ, G.: Leçons sur la Théorie Mathématique de l'Élasticité des Corps Solides. Paris: Bachelier.
- 1863 E 1. CLEBSCH, A.: Über die Reflexion an einer Kugelfläche. J. Reine Angew. Math. **61**, 195—262.
- 1882 E 1. CERRUTI, V.: Ricerche intorno all' equilibrio de' corpi elastici isotropi. Atti accad. Lincei Memorie (3) **13**, 81—123.
- 1885 E 1. BOUSSINESQ, M.: Application des potentiels à l'étude de l'équilibre et des mouvements des solides élastiques. Paris: Gauthier-Villars. See §§ II, VII.
- 1887 E 1. IBBETSON, W.: On the AIRY-MAXWELL solution of the equations of equilibrium of an isotropic elastic solid, under conservative forces. Proc. Lond. Math. Soc. **17**, 296—309 (1885—1886).
- 1889 E 1. JAERISCH, P.: Allgemeine Integration der Elasticitätsgleichungen für die Schwingungen und das Gleichgewicht isotroper Rotationskörper. J. Reine Angew. Math. **104**, 177—210.
- E 2. JAERISCH, P.: Zur Theorie der Elasticität isotroper Rotationskörper. Mitteil. Hamburger Math. Ges. **1**, 278—289.
- 1898 E 1. DUHEM, P.: Sur l'intégrale des équations des petits mouvements d'un solide isotrope. Mem. Soc. Sci. Bordeaux (5) **3**, 317—329.
- 1900 E 1. DUHEM, P.: Sur la généralisation d'un théorème de CLEBSCH. J. Math. Pures Appl. (5) **6**, 215—259.
- 1904 E 1. Lord KELVIN: Baltimore Lectures on Molecular Dynamics and the Wave Theory of Light (1884). London: C. J. Clay and Sons. See Lecture IV, p. 41 *et seq.*
- E 2. MARCOLONGO, R.: Teoria Matematica dello Equilibrio dei Corpi Elastici. Milano. See p. 236 *et seq.*
- E 3. TEDONE, O.: Saggio di una teoria generale delle equazioni dell' equilibrio elastico per un corpo isotropo. Annali di mat. (3) **10**, 13—64.
- 1914 E 1. GWYTHER, R.: The specification of the elements of stress. Part III. The definition of the dynamical specification and a test of the elastic specification. A chapter in elasticity. Mem. Manchester lit. philos. soc. **58**, No. 5 (21 pp.).

- 1915 E 1. KORN, A.: Über die beiden bisher zur Lösung der ersten Randwertaufgabe der Elastizitätstheorie eingeschlagenen Wege. *Annales acad. poly. Porto* **10**, 1—28.
- 1918 E 1. GWYTHYR, R. F.: The specification of stress, Part V. *Mem. Manchester lit. phil. soc.* **62**, No. 1, 11 pp. (1917—1918).
- 1919 E 1. SERINI, R.: Deformazione simmetriche dei corpi elastici. *Atti accad. Lincei* (5) **28**, 343—347.
- 1923 E 1. BURGATTI, P.: Sopra una soluzione molto generale dell' equazioni dell' equilibrio elastico. *Rend. accad. Bologna* (2) **27**, 66—73 (1922—1923).
- 1924 E 1. TIMPE, A.: Achsensymmetrische Deformation von Umdrehungskörpern. *Z. angew. Math. Mech.* **4**, 361—376.
- 1925 E 1. WEBER, C.: Achsensymmetrische Deformation von Umdrehungskörpern. *Z. angew. Math. Mech.* **5**, 466—468.
- 1926 E 1. BURGATTI, P.: Sopra due utili forme dell' integrale generale dell' equazioni per l'equilibrio dei solidi elastici isotrope. *Mem. accad. Bologna* (8) **3**, 63—67.
- 1928 E 1. TREFFTZ, E.: Mathematische Elastizitätstheorie. *Handbuch der Physik* **6**. Berlin: Springer. See p. 91 *et seqq.*
- E 2. LEVI-CIVITA, T.: *Fondamenti di Meccanica Relativistica*. Bologna: Zanichelli. vii + 185 pp. See §§ 25—26.
- 1930 E 1. GALERKIN, B.: Contribution à la solution générale du problème de la théorie de l'élasticité dans le cas de trois dimensions. *C. r. acad. sci. Paris* **190**, 1047—1048.
- E 2. GALERKIN, B.: On an investigation of stresses and deformations in elastic isotropic solids [in Russian]. *Doklad. Akad. Nauk SSSR* **1930 A**, 353—358.
- E 3. MALKIN, I.: Über einige neuere Arbeiten auf dem Gebiete der Elastizitätslehre. *Z. angew. Math. Mech.* **10**, 182—197.
- 1931 E 1. GALERKIN, B.: Sur l'équilibre élastique d'un plaque rectangulaire épaisse. *C. r. acad. sci. Paris* **193**, 563—571.
- E 2. GALERKIN, B.: Elastic rectangular and triangular thick plates with free edges, subjected to bending [in Russian]. *Doklad. Akad. Nauk SSSR* **1931**, 273—280.
- E 3. GALERKIN, B.: On the general solution of a problem in the theory of elasticity in three dimensions by means of stress and displacement functions [in Russian]. *Doklad. Akad. Nauk SSSR* **1931**, 281—286.
- 1932 E 1. GALERKIN, B.: On the investigation of stresses and deformation in a thick rectangular plate [in Russian]. *Isv. Nauchno-Issled. Inst. Gidrotekh.* **6**, 28—38.
- E 2. GALERKIN, B.: General solution of a problem on stresses and deformation in a thick circular plate and in a plate having the form of a circular sector [in Russian]. *Isv. Nauchno-Issled. Inst. Gidrotekh.* **7**, 1—6.
- E 3. PAPCOVITCH, P.: Solution générale des équations différentielles fondamentales d'élasticité, exprimée par trois fonctions harmoniques. *C. r. acad. sci. Paris* **195**, 513—515.
- E 4. PAPCOVITCH, P.: Expressions générales des composantes des tensions, ne renfermant comme fonctions arbitraires que des fonctions harmoniques. *C. r. acad. sci. Paris* **195**, 754—756.
- E 5. PAPKOVICH, P.: The representation of the general integral of the fundamental equations of the theory of elasticity in terms of harmonic functions [in Russian]. *Izv. Akad. Nauk SSSR, Fiz.-Mat. Ser.* **10**, 1425—1435.
- 1933 E 1. MARGUERRE, K.: Ebenes und achsensymmetrisches Problem der Elastizitätstheorie. *Z. angew. Math. Mech.* **13**, 437—438.
- 1934 E 1. BIEZENO, C.: Über die MARGUERRESche Spannungsfunktion. *Ing.-Archiv* **5**, 120—124.
- E 2. GALERKIN, B.: Contribution to the theory of an elastic cylindrical shell. *C. r. Doklady URSS* (2) **4**, 270—275.
- E 3. NEUBER, H.: Ein neuer Ansatz zur Lösung räumlicher Probleme der Elastizitätstheorie. Der Hohlkegel unter Einzellast als Beispiel. *Z. angew. Math. Mech.* **14**, 203—212.

- E 4. SOBRERO, L.: Nuovo metodo per lo studio dei problemi di elasticità con applicazione al problema della piastra forata. *Ricerche di Ingegneria* **2**, 255—264.
- 1935 E 1. SOBRERO, L.: Delle funzioni analoghe al potenziale intervenienti nella fisica-matematica. *Rend. Lincei*. (6) **21**, 448—454.
- E 2. STEUERMANN, E.: Sur une transformation des équations de la théorie d'élasticité. *J. inst. math. acad. sci. Ukraine* **1935**, Nos. 3—4, 35—40.
- E 3. WESTERGAARD, H.: General solution of the problem of elastostatics of an n -dimensional homogeneous isotropic solid in an n -dimensional space. *Bull. Amer. math. soc.* **41**, 695—699.
- 1936 E 1. MINDLIN, R.: Note on the GALERKIN and PAPKOVICH stress functions. *Bull. Amer. math. soc.* **42**, 373—376.
- E 2. ZANABONI, O.: Il problema della funzione delle tensioni in un sistema spaziale isotropo. *Bull. un. mat. Ital.* **15**, 71—76.
- 1937 E 1. NEUBER, H.: *Kerbspannungslehre*. Berlin: Springer. VII + 160 pp. See § 6.
- E 2. TOLOTTI, C.: Sui problemi di elasticità piana a funzione di AIRY polidroma. *Rend. Lincei* (6) **25**, 226—230.
- 1938 E 1. SLOBODIANSKI, M.: Stress functions for spatial problems in the theory of elasticity [in Russian]. *Uchenie Zap. Moskov. Gos. Univ.* **24**, 181—190.
- E 2. SLOBODIANSKI, M.: Expression of the solution of the differential equations of elasticity by means of one, two, and three functions and proof of the generality of these solutions [in Russian]. *Uchenie Zap. Moskov. Gos. Univ.* **24**, 191—202.
- 1939 E 1. BIEZENO, C. B., & R. GRAMMEL: *Technische Dynamik*. Berlin: Springer. See Ch. II, § 3.
- 1940 E 1. LEHKNITSKY, S.: Symmetric deformation and torsion of bodies of revolution with anisotropy of a special kind [in Russian]. *Priklad. math. mekh.* (3) **4**, No. 3, 43—60.
- E 2. LOCATELLI, P.: Sulla congruenza delle deformazioni. *Rend. ist. Lombardo* **73**, (4) 3, 457—464 (1939—1940).
- 1941 E 1. PLATRIER, C.: Sur l'intégration des équations indéfinies de l'équilibre élastique. *C. r. acad. sci. Paris* **212**, 749—751.
- 1942 E 1. SOUTHWELL, R.: Some practically important stress-systems in solids of revolution. *Proc. r. soc. London A* **180**, 367—396.
- 1944 E 1. CARRIER, G. F.: The thermal-stress and body-force problems for the infinite orthotropic solid. *Q. appl. math.* **2**, 31—36.
- E 2. FÖPPL, A. & L.: *Drang und Zwang 2*. Berlin: Oldenbourg. See p. 207 *et sqq.*
- 1946 E 1. BUTTY, E.: *Tratado de elasticidad teorico-tecnica 1*, Buenos Aires: Centro Estudiantes de Ingenieria. See Chs. VII—VIII.
- 1947 E 1. GRIOLI, G.: Struttura della funzione di AIRY nei sistemi molteplicemente connessi. *Giorn. mat* **77**, 119—144; trans., The structure of AIRY's stress function in multiply connected regions, NACA TM 1290, 1951.
- E 2. GUTMAN, C. G.: General solution of a problem in the theory of elasticity in generalized cylindrical coordinates [in Russian]. *Doklad. Akad. Nauk SSSR* **58**, 993—996.
- E 3. SADOWSKY, M., & E. STERNBERG: Stress concentration around an ellipsoidal cavity in an infinite body under arbitrary plane stress perpendicular to the axis of revolution of cavity. *J. appl. mech.* **69**, A 191—A 201.
- E 4. SHAPIRO, G.: Les fonctions des tensions dans un système arbitraire de coordonnées curvilignes. *C. r. Doklady acad. sci. URSS* **55**, 693—695.
- 1948 E 1. AYMERICH, G.: Trasformazione conforme delle funzioni biarmoniche ed applicazione all' elasticità piana. *Rend. sem. fac. sci. univ. Cagliari* **17**, 1—12 (1947).
- E 2. ELLIOT, H.: Three-dimensional stress distributions in hexagonal aeolotropic crystals. *Proc. Cambridge phil. soc.* **44**, 522—533.
- E 3. MOISIL, GR. C.: Sur une généralisation de la fonction d'AIRY. *Bull. éc. poly. Jassy* **3**, 156.

- E 4. TIMPE, A.: Torsionsfreie achsensymmetrische Deformation von Umdrehungskörpern und ihre Inversion. *Z. angew. Math. Mech.* **28**, 161—166.
- 1949 E 1. AGUARO, G.: Sul calcolo delle deformazioni di uno strato sferico elastico. *Rend. Lincei* (8) **7**, 289—297.
- E 2. FICHERA, G.: Sul calcolo delle deformazioni, dotate di simmetria assiale, di uno strato sferico elastico. *Rend. Lincei* (8) **6**, 582—590.
- E 3. FREIBERGER, W.: The uniform torsion of an incomplete tore. *Austral. J. sci. A* **2**, 354—375.
- E 4. FREIBERGER, W.: On the solution of the equilibrium equations of elasticity in general curvilinear coordinates. *Austral. J. sci. A* **2**, 483—492.
- E 5. IACOVACHE, M.: O extindere a metodei lui GALERKIN pentru sistemul ecuațiilor elasticității. *Acad. Române. Bul. Ști. A* **1**, 593—596.
- E 6. MARTIN, F.: Die Membran-Kugelschale unter Einzellasten. *Ing.-Arch.* **17**, 167—186.
- E 7. MOISIL, G.: Asupra sistemelor de ecuații cu derivate parțiale lineare și cu coeficienți constanți. *Acad. Române. Bul. Ști. A* **1**, 341—351.
- E 8. MOISIL, G.: Asupra formulelor lui GALERKIN în teoria elasticității. *Acad. Române. Bul. Ști. A* **1**, 587—592.
- E 9. MOISIL, G.: Un analog al vectorului lui GALERKIN în hidrodinamica lichidelor vâscoase. *Acad. Române. Bul. Ști. A* **1**, 803—812.
- E 10. STERNBERG, E., & M. SADOWSKY: Three-dimensional solution for the stress concentration around a circular hole in a plate of arbitrary thickness. *J. appl. mech.* **16**, 27—38.
- E 11. ZIEMBA, S.: Fonction des tensions aux coordonnées sphériques dans le cas d'une symétrie axiale des déformations et des tensions. *Archivum Mech. Stosow* **1**, 311—338.
- 1950 E 1. IONESCU-CAZIMIR, V.: Asupra ecuațiilor echilibrului termoelastic. I. Analogul vectorului lui GALERKIN. *Acad. Române. Bul. Ști. Ser. Mat. Fiz. Chim.* **2**, 589—595.
- E 2. LEHKNITSKY, S.: Theory of Elasticity of an Anisotropic Body [in Russian]. Moscow-Leningrad.
- E 3. MOISIL, A.: Asupra unui vector analog vectorului lui GALERKIN pentru echilibrul corpurilor elastice cu isotropie transversă. *Acad. Române. Bul. Ști. Ser. Mat. Fiz. Chim.* **2**, 207—210.
- E 4. TIMPE, A.: Spannungsfunktion für die von Kugel- und Kegelflächen begrenzten Körper und Kuppelproblem. *Z. angew. Math. Mech.* **30**, 50—61.
- 1951 E 1. IACOVACHE, M.: Relațiile între tensiuni într'un lichid vâscos incompresibil în mișcare lentă, permanentă. *Comun. acad. Române* **1**, 245—249.
- E 2. IONESCU-CAZIMIR, V.: Asupra ecuațiilor echilibrului termoelastic. II. Relațiile între tensiuni și temperatură. *Comun. acad. Române* **1**, 171—177.
- E 3. IONESCU-CAZIMIR, V.: Asupra ecuațiilor echilibrului termoelastic. III. Relațiile între tensiuni. *Comun. acad. Române* **1**, 385—390.
- E 4. KOÇO, P.: Asupra echilibrului unei clase de corpi visco-elastice. *Acad. Române. Bul. ști. sect. mat. fiz.* **3**, 221—243.
- E 5. STERNBERG, E., R. EUBANKS & M. SADOWSKY: On the stress-function approaches of BOUSSINESQ and TIMPE to the axisymmetric problem of elasticity theory. *J. appl. phys.* **22**, 1121—1124.
- E 6. TIMPE, A.: Spannungsfunktion achsensymmetrischer Deformationen in Zylinderkoordinaten. *Z. angew. Math. Mech.* **31**, 220—224.
- 1952 E 1. IACOVACHE, M.: Aplicarea funcțiilor monogene în sensul lui FEODOROV la teoria elasticității corpurilor cu izotropie transversă. *Rev. univ. „C. I. Parhon“ politeh. ser. ști. Mat.* **1**, No. 1, 58—60.
- E 2. IONESCU-CAZIMIR, V.: Asupra ecuațiilor echilibrului termoelastic plan. *Rev. univ. „C. I. Parhon“ politeh. ser. ști. Mat.* **1**, No. 1, 55—57.
- E 3. IONESCU-CAZIMIR, V.: Asupra ecuațiilor echilibrului termoelastic. IV. Cazul plan. *Acad. Române. Bul. ști. sect. mat. fiz.* **4**, 547—554.

- E 4. MALITA, M.: Asupra ecuațiilor de mișcare ale corpurilor elastice cu isotropie transversă. *Comun. acad. Române* **2**, 681—689.
- E 5. MOISIL, G. C.: Teoria preliminară a sistemelor de ecuații cu derivate parțiale cu coeficienți constanți. *Bul. științ. Acad. Pop. Române* **4**, 319ff.
- E 6. TEODORESCU, P.: Asupra teoriilor exacte a echilibrului suprafețelor cilindrice. *Acad. Române. Bul. ști. ser. mat. fiz.* **4**, 111—193.
- E 7. WESTERGAARD, H.: *Theory of Elasticity and Plasticity*. xii + 176 pp. Cambridge. See Ch. VI.
- 1953 E 1. ARZHANIKH, I.: Parametric representation of solutions of a system of linear functional equations in commutative operators [in Russian]. *Uspekhi Mat. Nauk* (2) **8**, No. 3 (55), 157—160.
- E 2. ARZHANIKH, I.: Studies on the mechanics of continuous media [in Russian]. *Akad. Nauk Uzbekh. SSR Trudi Inst. Mat. Mekh.* **9**, 60—101.
- E 3. BISHOP, R. E. D.: On dynamical problems of plane stress and plane strain. *Q. J. Mech. Appl. Math.* **6**, 250—254.
- E 4. BRDIČKA, M.: Equations of compatibility and stress functions in tensor form [in Russian]. *Czechosl. J. Phys.* **3**, 36—52.
- E 5. CHURIKOV, F.: On a form of the general solution of the equilibrium equations for the displacements in the theory of elasticity [in Russian]. *Prikl. Mat. Mekh.* **17**, 751—754.
- E 6. HU, H.: On the three-dimensional problems of the theory of elasticity of a transversely isotropic body. *Acta sci. Sinica* **2**, 145—151.
- E 7. SADOWSKY, M., & E. STERNBERG: Pure bending of an incomplete torus. *J. appl. mech.* **20**, 215—226.
- E 8. TRENIN, S.: On the solutions of the equilibrium equations of an axisymmetrical problem in the theory of elasticity [in Russian]. *Vestnik Moskov. Univ. Ser. Fiz.-Mat.* **8**, 7—13.
- 1954 E 1. BRDIČKA, M.: Covariant form of the general solution of the GALERKIN equations of elastic equilibrium [in Russian]. *Czechosl. J. Phys.* **4**, 246.
- E 2. EUBANKS, R., & E. STERNBERG: On the axisymmetric problem of elasticity theory for a medium with transverse isotropy. *J. Rational Mech. Anal.* **3**, 89—101.
- E 3. HU, H.: On the general theory of elasticity for a spherically isotropic medium. *Sci. Sinica* **3**, 247—260.
- E 4. IONESCU, D.: Asupra vectorului lui GALERKIN în teoria elasticității și în hidrodinamica fluidelor vâscoase. *Acad. Române. Bul. ști. sect. mat. fiz.* **6**, 555—571.
- E 5. KRÖNER, E.: Die Spannungsfunktionen der dreidimensionalen isotropen Elastizitätstheorie. *Z. Physik* **139**, 175—188. Correction, *Z. Physik* **143**, 175 (1955).
- E 6. LING, C.-B., & K.-L. YANG: On symmetrical strains in solids of revolution in curvilinear coordinates. *Ann. Acad. Sinica Taipei* **1**, 507—516.
- E 7. MINDLIN, R. D.: Force at a point in the interior of a semi-infinite solid. *Proc. 1st Midwest Conf. Solid Mech.* 1953, Univ. Illinois, 56—59.
- E 8. SLOBODIANSKI, M. G.: The general form of solutions of the equations of elasticity for simply connected and multiply connected domains, expressed by harmonic functions [in Russian]. *Prikl. Mat. Mekh.* **18**, 55—74.
- 1955 E 1. MARGUERRE, K.: Ansätze zur Lösung der Grundgleichungen der Elastizitätstheorie. *Z. angew. Math. Mech.* **35**, 241—263.
- 1956 E 1. BIOT, M. A.: Thermo-elasticity and irreversible thermodynamics. *J. Appl. Phys.* **27**, 240—253.
- E 2. EUBANKS, R. A., & E. STERNBERG: On the completeness of the BOUSSINESQ-PAPKOVICH stress functions. *J. Rational Mech. Anal.* **5**, 735—746.
- E 3. RADOK, J. R. M.: On the solution of problems of dynamic plane elasticity. *Q. Appl. Math.* **14**, 289—298.
- E 4. TEODORESCU, P. P.: On the plane problem of the elastodynamics. *Rev. Méc. Appl. Acad. Rep. Pop. Roumaine* **1**, 179ff.

- E 5. TEODORESCU, P.: On a general method of solving the plane problem of elastodynamics [in Romanian]. Com. Acad. R. P. Romine **6**, 795—801.
- 1957 E 1. NOLL, W.: Verschiebungsfunktionen für elastische Schwingungsprobleme. Z. angew. Math. Mech. **37**, 81—87.
- E 2. STERNBERG, E., & R. A. EUBANKS: On stress functions for elastokinetics and the integration of the repeated wave equation. Q. Appl. Math. **15**, 149—153.
- 1958 E 1. BLOKH, V. I.: On the representation of the general solution of the basic equations of the static theory of elasticity for an isotropic body with the aid of harmonic functions [in Russian]. Prikl. Mat. Mekh. **22**, 473—479. English, transl., PMM **22**, 659—668.
- E 2. DEEV, V. M.: The solution of the spatial problem in the theory of elasticity [in Ukrainian]. Dopov. Acad. Nauk Ukrain. RSR 1958, 29—32.
- E 3. DERESIEWICZ, H.: Solution of the equations of thermoelasticity. Proc. Third U.S. Nat. Cong. Appl. Mech.
- E 4. DUFFIN, R. J., & W. NOLL: On exterior boundary value problems in linear elasticity. Arch. Rational Mech. Anal. **2**, 191—196 (1958—1959).
- E 5. PREDELEANU, M.: Über die Verschiebungsfunktion für das achsensymmetrische Problem der Elastodynamik. Z. angew. Math. Mech. **38**, 402—405.
- 1960 E 1. STERNBERG, E.: On the integration of the equations of motion in the classical theory of elasticity. Arch. Rational Mech. Anal. (in press).

801 North College Avenue
Bloomington, Indiana

(Received July 1, 1959)

Infinitesimal Plane Strain in a Network of Elastic Cords

S. M. GENENSKY & R. S. RIVLIN

1. Introduction

In a previous paper^{*}, the problem has been considered of plane strain in a sheet of fabric consisting of a network formed by two families of parallel, inextensible, perfectly flexible cords, which cannot move relative to each other at their points of intersection. No restrictions were placed on the magnitude of the displacement gradients in the deformed sheet, except insofar as they are implied by the inextensibility of the cords.

In the present paper, we consider a similar problem in the case when the cords are extensible and elastic and the displacement gradients are sufficiently small so that terms involving the displacement gradients linearly may be neglected in comparison with unity, terms of second degree in the displacement gradients may be neglected in comparison with those of first degree and so on.

In Part I, the theory is developed without placing any limitations on the magnitude of the deformation undergone by the sheet and equations of equilibrium and expressions for the edge tractions are obtained for plane strain of the sheet. In Part II the relations obtained in Part I are simplified for the case when the displacement gradients are small compared with unity. It is shown (in § 8) how, in this case, the displacements throughout the sheet may be obtained when the edge displacements are specified. In § 9 it is shown how the displacements throughout the sheet may be obtained when the edge tractions are specified. Finally, in § 10 it is shown how the displacements throughout the sheet may be obtained in certain problems in which the boundary conditions are of the mixed type, the edge displacements being specified over part of the boundary and the edge tractions over the remainder of the boundary.

I. General Theory

2. Kinematic Considerations

The plane sheet of fabric described in § 1 is considered to undergo a plane deformation. We may describe the deformation with reference to a fixed rectangular Cartesian coordinate system x , the axes x_1 and x_2 of which are parallel to the bisectors of the angles between the cords in the undeformed state of the sheet, and are chosen so that the two families of cords make angles $\pm\alpha$ with the axis x_1 in the undeformed state of the sheet.

^{*} RIVLIN, R. S.: *J. Rational Mech. Anal.* **4**, 951 (1955).

If X_i and x_i ($i = 1, 2$)^{*} are the coordinates in the system x of a generic particle of the sheet in the undeformed and deformed states respectively, then the deformation of the sheet is described by relations of the form

$$x_i = x_i(X_j), \quad (2.1)$$

where the indicated functional dependence is single-valued, possesses a unique inverse, and is continuous and differentiable as many times as may be required, except possibly at isolated points or on isolated lines. Let P and Q be two neighboring particles of the sheet and let dS and ds be the distances between them in the undeformed and deformed states respectively. Then,

$$(dS^2) = dX_i dX_i \quad (2.2)$$

and

$$(ds)^2 = dx_i dx_i. \quad (2.3)$$

From (2.1), it follows that

$$dx_i = \frac{\partial x_i}{\partial X_j} dX_j. \quad (2.4)$$

Using (2.4), (2.3) becomes

$$(ds)^2 = \frac{\partial x_i}{\partial X_j} \frac{\partial x_i}{\partial X_k} dX_j dX_k. \quad (2.5)$$

Let L_i be the direction-cosines, in the rectangular Cartesian coordinate system x , of the line element PQ in the undeformed state of the sheet. Then,

$$L_i = dX_i/dS. \quad (2.6)$$

With (2.5), we obtain

$$\left(\frac{ds}{dS}\right)^2 = \frac{\partial x_i}{\partial X_j} \frac{\partial x_i}{\partial X_k} L_j L_k. \quad (2.7)$$

Let e_1 and e_2 be the fractional extensions, in the deformed state of the sheet, of cords which make angles α and $-\alpha$, respectively with the x_1 axis in the undeformed state. Then, taking $L_j = \cos \alpha, \sin \alpha$ in (2.7), we obtain

$$(1 + e_1)^2 = \left(\frac{\partial x_i}{\partial X_1} \cos \alpha + \frac{\partial x_i}{\partial X_2} \sin \alpha\right) \left(\frac{\partial x_i}{\partial X_1} \cos \alpha + \frac{\partial x_i}{\partial X_2} \sin \alpha\right). \quad (2.8)$$

Similarly, taking $L_j = \cos \alpha, -\sin \alpha$ in (2.7), we obtain

$$(1 + e_2)^2 = \left(\frac{\partial x_i}{\partial X_1} \cos \alpha - \frac{\partial x_i}{\partial X_2} \sin \alpha\right) \left(\frac{\partial x_i}{\partial X_1} \cos \alpha - \frac{\partial x_i}{\partial X_2} \sin \alpha\right). \quad (2.9)$$

Let Z_i denote the coordinates, in an oblique coordinate system, the axes of which are parallel to the directions of the two families of cords in the undeformed sheet, of a point which is at X_i in the system x . We then have

$$X_1 = (Z_1 + Z_2) \cos \alpha \quad \text{and} \quad X_2 = (Z_1 - Z_2) \sin \alpha. \quad (2.10)$$

From (2.10), we obtain

$$\begin{aligned} \frac{\partial}{\partial Z_1} &= \cos \alpha \frac{\partial}{\partial X_1} + \sin \alpha \frac{\partial}{\partial X_2} \\ \text{and} \quad \frac{\partial}{\partial Z_2} &= \cos \alpha \frac{\partial}{\partial X_1} - \sin \alpha \frac{\partial}{\partial X_2}. \end{aligned} \quad (2.11)$$

^{*} Latin subscripts will be considered to take the values 1, 2, throughout this paper.

Introducing (2.11) into (2.8) and (2.9), we obtain

$$(1 + e_1)^2 = \frac{\partial x_i}{\partial Z_1} \frac{\partial x_i}{\partial Z_1} \quad \text{and} \quad (1 + e_2)^2 = \frac{\partial x_i}{\partial Z_2} \frac{\partial x_i}{\partial Z_2}. \quad (2.12)$$

Let l_i be the direction-cosines of the line element PQ in the deformed state of the sheet, in the coordinate system x . Then

$$l_i = dx_i/ds. \quad (2.13)$$

With equations (2.4) and (2.6), we obtain

$$l_i = \frac{dS}{ds} \frac{\partial x_i}{\partial X_j} L_j. \quad (2.14)$$

3. Stresses

We consider next a sheet of fabric subjected to plane strain by the application of edge and surface tractions. The components of stress, t_{ij} , resulting from this deformation, are defined in the following manner; t_{i1} and t_{i2} are the components of the force per unit length in the positive directions of the x_1 and x_2 axes respectively, measured in the deformed state, exerted across an element of length at (x_1, x_2) normal to the x_i axis, by the material on the positive side of the element upon the material on the negative side of the element. From the continuity of the fabric and the equilibrium of moments, it follows that

$$t_{ij} = t_{ji}. \quad (3.1)$$

Let T_i be the components in the coordinate directions of the stress vector, measured in the deformed state, acting on an element of length at (x_1, x_2) with unit normal which has direction-cosines n_i in the system x . Then,

$$T_i = t_{ij} n_j. \quad (3.2)$$

4. The Relation Between Stress and Deformation

We assume that the sheet of fabric is deformed in accordance with equation (2.1) by forces acting in the plane of the sheet. Let τ_1 and τ_2 be the tensions, at a generic point P of the sheet, in the cords inclined initially at angles α and $-\alpha$ to the x_1 axis. Let β_1 and $-\beta_2$ be the inclinations of these cords to the x_1 axis in the deformed state.

Since linear elements of the cords of the two families, which have lengths dZ_1 and dZ_2 respectively in the undeformed state, have lengths $(1 + e_1) dZ_1$ and $(1 + e_2) dZ_2$ respectively in the deformed state, we have,

$$\begin{aligned} \cos \beta_1 &= \frac{1}{1+e_1} \frac{\partial x_1}{\partial Z_1}, & \sin \beta_1 &= \frac{1}{1+e_1} \frac{\partial x_2}{\partial Z_1}, \\ \cos \beta_2 &= \frac{1}{1+e_2} \frac{\partial x_1}{\partial Z_2}, & \sin \beta_2 &= -\frac{1}{1+e_2} \frac{\partial x_2}{\partial Z_2}. \end{aligned} \quad (4.1)$$

These results can be obtained by taking $L_1, L_2 = \cos \alpha, \sin \alpha$ and $l_1, l_2 = \cos \beta_1, \sin \beta_1$ and by taking $L_1, L_2 = \cos \alpha, -\sin \alpha$ and $l_1, l_2 = \cos \beta_2, -\sin \beta_2$ in equations (2.14) and employing the relations (2.11).

We denote by \bar{d} the intercept formed by adjacent cords of one family on a cord of the other family in the undeformed state. Then the intercept formed, in the deformed state, by adjacent cords at P of the family initially parallel to the Z_2 axis, on the cord through P , initially parallel to the Z_1 axis, is $\bar{d}(1+e_1)$. The number of cords of the Z_2 family intersecting the cord of the Z_1 family per unit length in the deformed state is therefore $1/\bar{d}(1+e_1)$. The stress vector on the Z_1 cord per unit length is thus $\tau_2/\bar{d}(1+e_1)$ in a direction parallel to the cords of the Z_2 family, *i.e.* in a direction having direction-cosines $\cos \beta_2$, $-\sin \beta_2$ in the coordinate system x . The components of the stress vector in the direction of the x_1 and x_2 axes are therefore

$$\frac{\tau_2 \cos \beta_2}{\bar{d}(1+e_1)} \quad \text{and} \quad -\frac{\tau_2 \sin \beta_2}{\bar{d}(1+e_1)} \quad (4.2)$$

respectively. The normal to the element of the Z_1 cord considered has direction-cosines $(\sin \beta_1, -\cos \beta_1)$. From equations (3.2) we see that the stress vector acting on such an element has components

$$t_{11} \sin \beta_1 - t_{12} \cos \beta_1 \quad \text{and} \quad t_{12} \sin \beta_1 - t_{22} \cos \beta_1 \quad (4.3)$$

in the x_1 and x_2 directions respectively. We thus have from (4.2) and (4.3)

$$\begin{aligned} t_{11} \sin \beta_1 - t_{12} \cos \beta_1 &= \frac{\tau_2 \cos \beta_2}{\bar{d}(1+e_1)} \\ \text{and} \quad t_{12} \sin \beta_1 - t_{22} \cos \beta_1 &= -\frac{\tau_2 \sin \beta_2}{\bar{d}(1+e_1)}. \end{aligned} \quad (4.4)$$

In a similar fashion, by considering the stress vector acting on an element of the Z_2 cord through P , we obtain

$$\begin{aligned} t_{11} \sin \beta_2 + t_{12} \cos \beta_2 &= \frac{\tau_1 \cos \beta_1}{\bar{d}(1+e_2)} \\ \text{and} \quad t_{12} \sin \beta_2 + t_{22} \cos \beta_2 &= \frac{\tau_1 \sin \beta_1}{\bar{d}(1+e_2)}. \end{aligned} \quad (4.5)$$

Solving equations (4.4) and (4.5) for t_{11} , t_{22} and t_{12} , we obtain

$$\begin{aligned} t_{11} &= \frac{1}{\bar{d} \sin(\beta_1 + \beta_2)} \left(\frac{\tau_1 \cos^2 \beta_1}{1+e_2} + \frac{\tau_2 \cos^2 \beta_2}{1+e_1} \right), \\ t_{22} &= \frac{1}{\bar{d} \sin(\beta_1 + \beta_2)} \left(\frac{\tau_1 \sin^2 \beta_1}{1+e_2} + \frac{\tau_2 \sin^2 \beta_2}{1+e_1} \right) \\ \text{and} \quad t_{12} &= \frac{1}{\bar{d} \sin(\beta_1 + \beta_2)} \left(\frac{\tau_1 \cos \beta_1 \sin \beta_1}{1+e_2} - \frac{\tau_2 \cos \beta_2 \sin \beta_2}{1+e_1} \right). \end{aligned} \quad (4.6)$$

Let d be the distance between adjacent cords of each family in the undeformed state of the sheet. Then,

$$\bar{d} = d/\sin 2\alpha. \quad (4.7)$$

Substituting from (4.1) and (4.7) in (4.6), we obtain

$$\begin{aligned} t_{11} &= \frac{\sin 2\alpha}{dA} \left[\frac{\tau_1}{1+e_1} \left(\frac{\partial x_1}{\partial Z_1} \right)^2 + \frac{\tau_2}{1+e_2} \left(\frac{\partial x_1}{\partial Z_2} \right)^2 \right], \\ t_{22} &= \frac{\sin 2\alpha}{dA} \left[\frac{\tau_1}{1+e_1} \left(\frac{\partial x_2}{\partial Z_1} \right)^2 + \frac{\tau_2}{1+e_2} \left(\frac{\partial x_2}{\partial Z_2} \right)^2 \right] \\ \text{and} \quad t_{12} &= \frac{\sin 2\alpha}{dA} \left[\frac{\tau_1}{1+e_1} \frac{\partial x_1}{\partial Z_1} \frac{\partial x_2}{\partial Z_1} + \frac{\tau_2}{1+e_2} \frac{\partial x_1}{\partial Z_2} \frac{\partial x_2}{\partial Z_2} \right], \end{aligned} \quad (4.8)$$

where Δ is defined by

$$\Delta = \frac{\partial x_1}{\partial Z_2} \frac{\partial x_2}{\partial Z_1} - \frac{\partial x_1}{\partial Z_1} \frac{\partial x_2}{\partial Z_2}. \quad (4.9)$$

5. The Equations of Motion

The equations of motion for the sheet are

$$\frac{\partial t_{ij}}{\partial x_j} + F_i = \rho \frac{\partial^2 x_i}{\partial t^2}, \quad (5.1)$$

where F_i denotes the applied force per unit area, measured in the deformed state, and ρ denotes the mass of fabric per unit area, measured in the deformed state. We shall limit our analysis to static problems, so that $\partial^2 x_i / \partial t^2 = 0$ and equation (5.1) can then be written as

$$\frac{\partial t_{ij}}{\partial Z_k} \frac{\partial Z_k}{\partial x_j} + F_i = 0. \quad (5.2)$$

We have

$$\frac{\partial x_i}{\partial Z_k} \frac{\partial Z_k}{\partial x_j} = \delta_{ij}. \quad (5.3)$$

Regarding (5.3) as an equation for the determination of $\partial Z_k / \partial x_j$, we obtain the solution as

$$\frac{\partial Z_k}{\partial x_j} = \frac{1}{\Delta} \left(\frac{\partial x_k}{\partial Z_j} - \delta_{jk} \frac{\partial x_p}{\partial Z_p} \right), \quad (5.4)$$

where Δ is defined by equation (4.9).

Introducing (5.4) into (5.2), we obtain

$$\frac{\partial t_{ij}}{\partial Z_k} \left(\frac{\partial x_k}{\partial Z_j} - \delta_{jk} \frac{\partial x_p}{\partial Z_p} \right) + \Delta F_i = 0. \quad (5.5)$$

The expressions (4.8) may be introduced into (5.5) and the resulting equations simplified to yield

$$\begin{aligned} \text{and} \quad & \frac{\partial}{\partial Z_1} \left(\frac{\tau_1}{1+e_1} \frac{\partial x_1}{\partial Z_1} \right) + \frac{\partial}{\partial Z_2} \left(\frac{\tau_2}{1+e_2} \frac{\partial x_1}{\partial Z_2} \right) - \frac{d \Delta F_1}{\sin 2\alpha} = 0 \\ & \frac{\partial}{\partial Z_1} \left(\frac{\tau_1}{1+e_1} \frac{\partial x_2}{\partial Z_1} \right) + \frac{\partial}{\partial Z_2} \left(\frac{\tau_2}{1+e_2} \frac{\partial x_2}{\partial Z_2} \right) - \frac{d \Delta F_2}{\sin 2\alpha} = 0. \end{aligned} \quad (5.6)$$

If no surface forces act on the sheet, *i.e.* $F_1 = F_2 = 0$, equations (5.6) become

$$\begin{aligned} \text{and} \quad & \frac{\partial}{\partial Z_1} \left(\frac{\tau_1}{1+e_1} \frac{\partial x_1}{\partial Z_1} \right) + \frac{\partial}{\partial Z_2} \left(\frac{\tau_2}{1+e_2} \frac{\partial x_1}{\partial Z_2} \right) = 0 \\ & \frac{\partial}{\partial Z_1} \left(\frac{\tau_1}{1+e_1} \frac{\partial x_2}{\partial Z_1} \right) + \frac{\partial}{\partial Z_2} \left(\frac{\tau_2}{1+e_2} \frac{\partial x_2}{\partial Z_2} \right) = 0. \end{aligned} \quad (5.7)$$

6. The Edge Traction

We now consider the plane sheet of fabric to be bounded by a simply-connected closed contour C , which is everywhere continuous and has a continuously-turning tangent except possibly at a finite number of points. We consider that the sheet is subjected to edge tractions in its plane, applied to the contour C .

Let Q be some reference particle and P a generic particle of the boundary C . We denote by Σ and σ the distance from Q to P , measured along C in the counter-clockwise sense, in the undeformed and deformed states of the sheet respectively. We further denote by N_i and n_i the direction-cosines of the outward-drawn normal to C at P in the undeformed and deformed states of the sheet respectively, and by M_i and m_i the direction-cosines of the tangent to C at P in these states. It follows immediately that

$$N_1 = M_2, \quad N_2 = -M_1 \quad (6.1)$$

and

$$n_1 = m_2, \quad n_2 = -m_1. \quad (6.2)$$

Let the components of the edge tractions, per unit length of contour, measured in the undeformed and deformed states of the sheet, be ξ_i and $\bar{\xi}_i$ respectively. Then,

$$\xi_i d\Sigma = \bar{\xi}_i d\sigma, \quad (6.3)$$

since $d\Sigma$ and $d\sigma$ are corresponding elements of length on the contour in the undeformed and deformed states respectively.

The components of edge traction $\bar{\xi}_i$ are given in terms of the stress components t_{ij} by

$$\bar{\xi}_i = t_{ij} n_j. \quad (6.4)$$

With (6.2), (6.4) yields

$$\bar{\xi}_1 = t_{11} m_2 - t_{12} m_1 \quad \text{and} \quad \bar{\xi}_2 = t_{12} m_2 - t_{22} m_1. \quad (6.5)$$

From (2.14), we obtain

$$m_i = \frac{d\Sigma}{d\sigma} \frac{\partial x_i}{\partial X_j} M_j. \quad (6.6)$$

Introducing (6.6) and (6.3) into (6.5), we obtain

$$\begin{aligned} \xi_1 &= \left(\frac{\partial x_2}{\partial X_j} t_{11} - \frac{\partial x_1}{\partial X_j} t_{12} \right) M_j \\ \text{and} \quad \xi_2 &= \left(\frac{\partial x_2}{\partial X_j} t_{12} - \frac{\partial x_1}{\partial X_j} t_{22} \right) M_j. \end{aligned} \quad (6.7)$$

Introducing the expressions (4.8) for t_{ij} into (6.7), we obtain

$$\begin{aligned} \xi_1 &= \frac{\sin 2\alpha}{d\Delta} \left[\frac{\tau_1}{1+e_1} \frac{\partial x_1}{\partial Z_1} \left(\frac{\partial x_1}{\partial Z_1} \frac{\partial x_2}{\partial X_j} - \frac{\partial x_2}{\partial Z_1} \frac{\partial x_1}{\partial X_j} \right) + \right. \\ &\quad \left. + \frac{\tau_2}{1+e_2} \frac{\partial x_1}{\partial Z_2} \left(\frac{\partial x_1}{\partial Z_2} \frac{\partial x_2}{\partial X_j} - \frac{\partial x_2}{\partial Z_2} \frac{\partial x_1}{\partial X_j} \right) \right] M_j \\ \text{and} \quad \xi_2 &= \frac{\sin 2\alpha}{d\Delta} \left[\frac{\tau_1}{1+e_1} \frac{\partial x_2}{\partial Z_1} \left(\frac{\partial x_1}{\partial Z_1} \frac{\partial x_2}{\partial X_j} - \frac{\partial x_2}{\partial Z_1} \frac{\partial x_1}{\partial X_j} \right) + \right. \\ &\quad \left. + \frac{\tau_2}{1+e_2} \frac{\partial x_2}{\partial Z_2} \left(\frac{\partial x_1}{\partial Z_2} \frac{\partial x_2}{\partial X_j} - \frac{\partial x_2}{\partial Z_2} \frac{\partial x_1}{\partial X_j} \right) \right] M_j. \end{aligned} \quad (6.8)$$

With (2.11) and (4.9), equations (6.8) yield

$$\begin{aligned} \xi_1 &= \frac{1}{d} \left[\frac{\tau_1}{1+e_1} \frac{\partial x_1}{\partial Z_1} (M_2 \cos \alpha - M_1 \sin \alpha) + \frac{\tau_2}{1+e_2} \frac{\partial x_1}{\partial Z_2} (M_2 \cos \alpha + M_1 \sin \alpha) \right] \\ \text{and} \quad \xi_2 &= \frac{1}{d} \left[\frac{\tau_1}{1+e_1} \frac{\partial x_2}{\partial Z_1} (M_2 \cos \alpha - M_1 \sin \alpha) + \frac{\tau_2}{1+e_2} \frac{\partial x_2}{\partial Z_2} (M_2 \cos \alpha + M_1 \sin \alpha) \right]. \end{aligned} \quad (6.9)$$

II. Infinitesimal Deformations

7. Development of Basic Equations

Let u_i be the components in the coordinate system x of the displacement undergone in the deformation by a generic particle of the sheet. Then, we have

$$x_i = X_i + u_i. \quad (7.1)$$

We shall assume that the deformation is sufficiently small so that terms of the second degree in the displacement gradients $\partial u_i / \partial X_j$ may be neglected in comparison with those of the first degree. From (2.11), it is apparent that we may neglect terms of the second degree in $\partial u_i / \partial Z_j$ in comparison with those of the first degree. It is also apparent that we may neglect terms of the second degree in e_i in comparison with those of first degree.

With this approximation, we can write equations (2.12) as

$$\begin{aligned} \cos \alpha \frac{\partial u_1}{\partial Z_1} + \sin \alpha \frac{\partial u_2}{\partial Z_1} &= e_1 \\ \cos \alpha \frac{\partial u_1}{\partial Z_2} - \sin \alpha \frac{\partial u_2}{\partial Z_2} &= e_2. \end{aligned} \quad (7.2)$$

We shall assume that the fractional extensions e_1 and e_2 in the cords of the two families are related to the tensions τ_1 and τ_2 in them by the formulae

$$\tau_1 = k e_1 \quad \text{and} \quad \tau_2 = k e_2, \quad (7.3)$$

where k is a constant.

Employing the relations (7.1) and (7.3) in equations (5.7) and neglecting terms of the second degree in e_1 , e_2 and $\partial u_i / \partial Z_j$, we obtain the equations of equilibrium, in the absence of surface forces as,

$$\begin{aligned} \frac{\partial e_1}{\partial Z_1} + \frac{\partial e_2}{\partial Z_2} &= 0 \\ \frac{\partial e_1}{\partial Z_1} - \frac{\partial e_2}{\partial Z_2} &= 0. \end{aligned} \quad (7.4)$$

From (7.4),

$$\frac{\partial e_1}{\partial Z_1} = \frac{\partial e_2}{\partial Z_2} = 0, \quad (7.5)$$

so that

$$e_1 = e_1(Z_2) \quad \text{and} \quad e_2 = e_2(Z_1); \quad (7.6)$$

i.e. the fractional extension—and hence the tension—is constant along each cord of the sheet, but may vary from cord to cord. Bearing this in mind, we now integrate equations (7.2) to obtain

$$\begin{aligned} u_1 \cos \alpha + u_2 \sin \alpha &= Z_1 e_1(Z_2) + p_1(Z_2) \\ u_1 \cos \alpha - u_2 \sin \alpha &= Z_2 e_2(Z_1) + p_2(Z_1), \end{aligned} \quad (7.7)$$

where $p_1(Z_2)$ and $p_2(Z_1)$ are arbitrary functions of their arguments. From (7.7), we obtain

$$\begin{aligned} 2u_1 \cos \alpha &= Z_1 e_1(Z_2) + Z_2 e_2(Z_1) + p_1(Z_2) + p_2(Z_1) \\ 2u_2 \sin \alpha &= Z_1 e_1(Z_2) - Z_2 e_2(Z_1) + p_1(Z_2) - p_2(Z_1). \end{aligned} \quad (7.8)$$

Again, employing the relations (7.1) and (7.3) in equations (4.8) and (4.9) and making the approximations employed above, we obtain the following expressions for the stress components t_{ij} :

$$\begin{aligned} t_{11} &= \frac{k \cos^2 \alpha}{d} (e_1 + e_2), & t_{22} &= \frac{k \sin^2 \alpha}{d} (e_1 + e_2) \\ \text{and} & & & \\ t_{12} &= t_{21} = \frac{k \cos \alpha \sin \alpha}{d} (e_1 - e_2). \end{aligned} \quad (7.9)$$

In a similar manner, we obtain, from equations (6.9), the following expressions for the edge tractions ξ_i :

$$\begin{aligned} \xi_1 &= \frac{k \cos \alpha}{d} [(e_1 + e_2) M_2 \cos \alpha - (e_1 - e_2) M_1 \sin \alpha] \\ \text{and} & & & \\ \xi_2 &= \frac{k \sin \alpha}{d} [(e_1 - e_2) M_2 \cos \alpha - (e_1 + e_2) M_1 \sin \alpha]. \end{aligned} \quad (7.10)$$

8. Solution of the Problem when Edge Displacements are Specified

We consider a sheet which in the undeformed state is bounded by a simply-connected closed contour C which is everywhere continuous and has a continuously-turning tangent, except possibly at a finite number of points, as shown in Fig. 1. Let P be an interior point of the sheet with coordinates (Z_1, Z_2) in the oblique system formed by two intersecting cords. Two cords pass through P cutting C at four distinct points Q, R, S and T , which have coordinates (a_1, Z_2) , (Z_1, a_2) , (b_1, Z_2) and (Z_1, b_2) respectively.

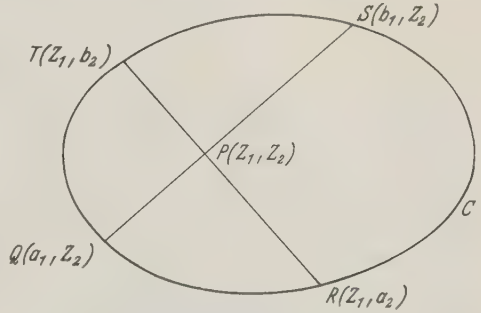


Fig. 1. Plane sheet bounded by simply-connected closed contour

From the first of equations (7.7), it follows that

$$\cos \alpha u_1(a_1, Z_2) + \sin \alpha u_2(a_1, Z_2) = a_1 e_1(Z_2) + p_1(Z_2) \text{ at } Q \quad (8.1)$$

and

$$\cos \alpha u_1(b_1, Z_2) + \sin \alpha u_2(b_1, Z_2) = b_1 e_1(Z_2) + p_1(Z_2) \text{ at } S. \quad (8.2)$$

Solving (8.1) and (8.2) for $e_1(Z_2)$ and $p_1(Z_2)$, we obtain

$$e_1(Z_2) = \frac{1}{a_1 - b_1} \{ \cos \alpha [u_1(a_1, Z_2) - u_1(b_1, Z_2)] + \sin \alpha [u_2(a_1, Z_2) - u_2(b_1, Z_2)] \} \quad (8.3)$$

and

$$p_1(Z_2) = \frac{1}{b_1 - a_1} \{ \cos \alpha [b_1 u_1(a_1, Z_2) - a_1 u_1(b_1, Z_2)] + \sin \alpha [b_1 u_2(a_1, Z_2) - a_1 u_2(b_1, Z_2)] \}.$$

Again, from the second of equations (7.7), it follows that

$$\cos \alpha u_1(Z_1, a_2) - \sin \alpha u_2(Z_1, a_2) = a_2 e_2(Z_1) + p_2(Z_1) \text{ at } R \quad (8.4)$$

and

$$\cos \alpha u_1(Z_1, b_2) - \sin \alpha u_2(Z_1, b_2) = b_2 e_2(Z_1) + p_2(Z_1) \text{ at } T. \quad (8.5)$$

Solving (8.4) and (8.5) for $e_2(Z_1)$ and $\phi_2(Z_1)$, we obtain

$$e_2(Z_1) = \frac{1}{a_2 - b_2} \{ \cos \alpha [u_1(Z_1, a_2) - u_1(Z_1, b_2)] - \sin \alpha [u_2(Z_1, a_2) - u_2(Z_1, b_2)] \} \quad (8.6)$$

and

$$\phi_2(Z_1) = \frac{1}{b_2 - a_2} \{ \cos \alpha [b_2 u_1(Z_1, a_2) - a_2 u_1(Z_1, b_2)] - \sin \alpha [b_2 u_2(Z_1, a_2) - a_2 u_2(Z_1, b_2)] \}.$$

Introducing (8.3) and (8.6) into (7.8), we arrive at

$$\begin{aligned} u_1(Z_1, Z_2) = & \frac{1}{2(a_1 - b_1)} \{ (Z_1 - b_1) [u_1(a_1, Z_2) + \tan \alpha u_2(a_1, Z_2)] - \\ & - (Z_1 - a_1) [u_1(b_1, Z_2) + \tan \alpha u_2(b_1, Z_2)] \} + \\ & + \frac{1}{2(a_2 - b_2)} \{ (Z_2 - b_2) [u_1(Z_1, a_2) - \tan \alpha u_2(Z_1, a_2)] - \\ & - (Z_2 - a_2) [u_1(Z_1, b_2) - \tan \alpha u_2(Z_1, b_2)] \} \end{aligned} \quad (8.7)$$

and

$$\begin{aligned} u_2(Z_1, Z_2) = & \frac{1}{2(a_1 - b_1)} \{ (Z_1 - b_1) [\cot \alpha u_1(a_1, Z_2) + u_2(a_1, Z_2)] - \\ & - (Z_1 - a_1) [\cot \alpha u_1(b_1, Z_2) + u_2(b_1, Z_2)] \} - \\ & - \frac{1}{2(a_2 - b_2)} \{ (Z_2 - b_2) [\cot \alpha u_1(Z_1, a_2) - u_2(Z_1, a_2)] - \\ & - (Z_2 - a_2) [\cot \alpha u_1(Z_1, b_2) - u_2(Z_1, b_2)] \}. \end{aligned}$$

Thus if u_i are known on the entire boundary C , they are uniquely determined throughout the sheet by equations (8.7).

We shall now suppose that a portion AB of the contour C is formed by a cord $Z_2 = \text{const.}$, as shown in Fig. 2(i). It follows from the first of equations (7.7) that $u_1 \cos \alpha + u_2 \sin \alpha$ is a linear function of Z_1 on AB . Consequently, if u_1 is specified on AB and u_2 is specified at two points of it, u_2 can be determined over the whole of AB . Again, if u_2 is specified on AB and u_1 is specified at two points of it, u_1 can be determined over the whole of AB . It follows that if u_1 is specified over the whole of the contour C and u_2 is specified on the whole of

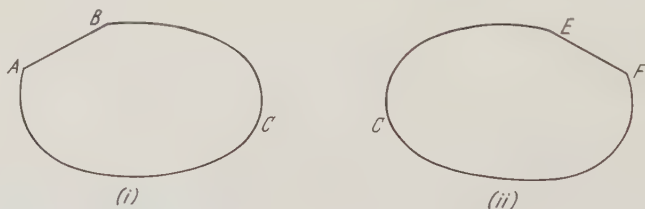


Fig. 2. Plane sheets bounded by simply-connected closed contour part of which is formed by a cord

the contour C except AB and at two points of AB (which may be the end-points A and B) then u_2 can be determined over the whole of the contour C and equations (8.7) used to determine u_1 and u_2 throughout the sheet. Alternatively, if u_2 is specified over the whole of the contour C and u_1 is specified over the whole of the contour C except AB and at two points of AB (which may be the end-points A and B) then u_1 can be determined over the whole of the contour C and equations (8.7) used to determine u_1 and u_2 throughout the sheet.

We now consider that a portion EF of the contour C is formed by a cord $Z_1 = \text{const}$, as shown in Fig. 2(ii). It follows from the second of equations (7.7) that $u_1 \cos \alpha - u_2 \sin \alpha$ is a linear function of Z_2 on EF . By an argument similar to that used in the previous case, we see that if one of the displacement components u_i is specified on the whole of the contour C and the other at all points of the contour C except EF and at two points of EF (which may be the end-points E and F), then both components of u_i can be determined over the whole of the contour C and equations (8.7) used to determine $u_i(Z_1, Z_2)$ throughout the sheet bounded by C .

It follows from this result that if the contour C is formed entirely by four cords, then if u_1 is specified at all points of C and u_2 at the vertices, u_2 may be determined over the whole of C . Alternatively, if u_2 is specified at all points of C and u_1 at the vertices, u_1 may be determined over the whole of C . In either case, equations (8.7) may then be used to determine u_1 and u_2 throughout the sheet.

9. Solution of the Problem when Edge Traction are Specified

We consider the sheet of material shown in Fig. 1 and described in § 8. We assume that the edge traction ξ_i are specified over the whole boundary C of the sheet. Then, on C we have, solving equations (7.10) for $e_1(Z_2)$ and $e_2(Z_1)$,

$$e_1(Z_2) = \frac{\bar{d}[\xi_1(Z_1, Z_2) \sin \alpha + \xi_2(Z_1, Z_2) \cos \alpha]}{k[M_2(Z_1, Z_2) \cos \alpha - M_1(Z_1, Z_2) \sin \alpha]} \quad (9.1)$$

and

$$e_2(Z_1) = \frac{\bar{d}[\xi_1(Z_1, Z_2) \sin \alpha - \xi_2(Z_1, Z_2) \cos \alpha]}{k[M_2(Z_1, Z_2) \cos \alpha + M_1(Z_1, Z_2) \sin \alpha]},$$

provided that

$$M_2(Z_1, Z_2) \cos \alpha \pm M_1(Z_1, Z_2) \sin \alpha \neq 0. \quad (9.2)$$

The condition (9.2) will, of course, be violated at points on C where the tangent to C is parallel to a cord direction.

If at a point on C the tangent to C is parallel to $Z_2 = \text{const}$, then $M_1, M_2 = \cos \alpha, \sin \alpha$. We then obtain, from (7.10) and (4.7),

$$e_2(Z_1) = \frac{\bar{d}}{k \cos \alpha} \xi_1(Z_1, Z_2) = -\frac{\bar{d}}{k \sin \alpha} \xi_2(Z_1, Z_2). \quad (9.3)$$

Again, if at a point on C the tangent to C is parallel to $Z_1 = \text{const}$, then $M_1, M_2 = \cos \alpha, -\sin \alpha$. We then obtain, from (7.10) and (4.7),

$$e_1(Z_2) = -\frac{\bar{d}}{k \cos \alpha} \xi_1(Z_1, Z_2) = -\frac{\bar{d}}{k \sin \alpha} \xi_2(Z_1, Z_2). \quad (9.4)$$

If the points Q, R, S and T are such that the condition (9.2) is not violated at any of them, then, from (9.1), we obtain

$$e_1(Z_2) = \frac{\bar{d}[\xi_1(a_1, Z_2) \sin \alpha + \xi_2(a_1, Z_2) \cos \alpha]}{k[M_2(a_1, Z_2) \cos \alpha - M_1(a_1, Z_2) \sin \alpha]} \quad \text{at } Q, \quad (9.5)$$

$$e_2(Z_1) = \frac{\bar{d}[\xi_1(Z_1, a_2) \sin \alpha - \xi_2(Z_1, a_2) \cos \alpha]}{k[M_2(Z_1, a_2) \cos \alpha + M_1(Z_1, a_2) \sin \alpha]} \quad \text{at } R, \quad (9.6)$$

$$e_1(Z_2) = \frac{\bar{d}[\xi_1(b_1, Z_2) \sin \alpha + \xi_2(b_1, Z_2) \cos \alpha]}{k[M_2(b_1, Z_2) \cos \alpha - M_1(b_1, Z_2) \sin \alpha]} \quad \text{at } S \quad (9.7)$$

and

$$e_2(Z_1) = \frac{\bar{d}[\xi_1(Z_1, b_2) \sin \alpha - \xi_2(Z_1, b_2) \cos \alpha]}{k[M_2(Z_1, b_2) \cos \alpha + M_1(Z_1, b_2) \sin \alpha]} \quad \text{at } T. \quad (9.8)$$

Comparing equations (9.5) and (9.7) and equations (9.6) and (9.8), we see that ξ_i cannot be arbitrarily specified on the boundary, but must satisfy the conditions

$$\text{and } \frac{\xi_1(a_1, Z_2) \sin \alpha + \xi_2(a_1, Z_2) \cos \alpha}{M_2(a_1, Z_2) \cos \alpha - M_1(a_1, Z_2) \sin \alpha} = \frac{\xi_1(b_1, Z_2) \sin \alpha + \xi_2(b_1, Z_2) \cos \alpha}{M_2(b_1, Z_2) \cos \alpha - M_1(b_1, Z_2) \sin \alpha} \quad (9.9)$$

$$\frac{\xi_1(Z_1, a_2) \sin \alpha - \xi_2(Z_1, a_2) \cos \alpha}{M_2(Z_1, a_2) \cos \alpha + M_1(Z_1, a_2) \sin \alpha} = \frac{\xi_1(Z_1, b_2) \sin \alpha - \xi_2(Z_1, b_2) \cos \alpha}{M_2(Z_1, b_2) \cos \alpha + M_1(Z_1, b_2) \sin \alpha}.$$

Introducing (9.5) and (9.6) into (7.8), we obtain

$$2 \cos \alpha u_1(Z_1, Z_2) = Z_1 \frac{\bar{d}[\xi_1(a_1, Z_2) \sin \alpha + \xi_2(a_1, Z_2) \cos \alpha]}{k [M_2(a_1, Z_2) \cos \alpha - M_1(a_1, Z_2) \sin \alpha]} +$$

$$+ Z_2 \frac{\bar{d}[\xi_1(Z_1, a_2) \sin \alpha - \xi_2(Z_1, a_2) \cos \alpha]}{k [M_2(Z_1, a_2) \cos \alpha + M_1(Z_1, a_2) \sin \alpha]} +$$

$$+ \phi_1(Z_2) + \phi_2(Z_1) \quad (9.10)$$

and

$$2 \sin \alpha u_2(Z_1, Z_2) = Z_1 \frac{\bar{d}[\xi_1(a_1, Z_2) \sin \alpha + \xi_2(a_1, Z_2) \cos \alpha]}{k [M_2(a_1, Z_2) \cos \alpha - M_1(a_1, Z_2) \sin \alpha]} -$$

$$- Z_2 \frac{\bar{d}[\xi_1(Z_1, a_2) \sin \alpha - \xi_2(Z_1, a_2) \cos \alpha]}{k [M_2(Z_1, a_2) \cos \alpha + M_1(Z_1, a_2) \sin \alpha]} +$$

$$+ \phi_1(Z_2) - \phi_2(Z_1).$$

We see that u_i are undetermined to the extent of a displacement field u'_i of the form

$$2 \cos \alpha u'_1 = \phi_1(Z_2) + \phi_2(Z_1) \quad (9.11)$$

$$2 \sin \alpha u'_2 = \phi_1(Z_2) - \phi_2(Z_1),$$

which is an arbitrary (small) displacement field for which cord lengths remain unchanged. This indeterminacy of the displacement field u_i is not surprising in view of the fact that the fabric may be subjected to any deformation in which the lengths of cord elements are unchanged, without forces being applied.

If at one or more of the points Q, R, S, T , the tangent to C is parallel to $Z_2 = \text{const}$, we obtain $e_2(Z_1)$ from (9.3) and a knowledge of either ξ_1 or ξ_2 at that point. Similarly, if at one or more of the points Q, R, S, T , the tangent to C is parallel to $Z_1 = \text{const}$, we obtain $e_1(Z_2)$ from (9.4) and a knowledge of either ξ_1 or ξ_2 at that point. These relations may then be employed together with relations chosen from (9.5) to (9.8), appropriate to those of the points Q, R, S, T at which the tangent to C is not parallel to a cord, to provide expressions for $e_1(Z_2)$ and $e_2(Z_1)$ at P . These expressions may then be introduced into (7.8) to give expressions for $u_1(Z_1, Z_2)$ and $u_2(Z_1, Z_2)$ at P . By way of illustration, we consider below two cases.

(i) *Sheet bounded by cords.* We consider the boundary C of the sheet to be formed by the cords $Z_2 = a_2, Z_2 = b_2, Z_1 = a_1, Z_1 = b_1$, as shown in Fig. 3. Let P be a generic point (Z_1, Z_2) of the sheet and let Q, R, S, T be the intersections

with C of the cords through P . We then have, from (9.3) and (9.4),

$$\begin{aligned}
 e_2(Z_1) &= \frac{\bar{d}}{k \cos \alpha} \xi_1(Z_1, a_2) = -\frac{\bar{d}}{k \sin \alpha} \xi_2(Z_1, a_2) \quad \text{at } T, \\
 e_2(Z_1) &= \frac{\bar{d}}{k \cos \alpha} \xi_1(Z_1, b_2) = -\frac{\bar{d}}{k \sin \alpha} \xi_2(Z_1, b_2) \quad \text{at } R, \\
 e_1(Z_2) &= -\frac{\bar{d}}{k \cos \alpha} \xi_1(a_1, Z_2) = -\frac{\bar{d}}{k \sin \alpha} \xi_2(a_1, Z_2) \quad \text{at } Q, \\
 e_1(Z_2) &= -\frac{\bar{d}}{k \cos \alpha} \xi_1(b_1, Z_2) = -\frac{\bar{d}}{k \sin \alpha} \xi_2(b_1, Z_2) \quad \text{at } S.
 \end{aligned} \tag{9.12}$$

Introducing (9.12) into (7.8), we obtain expressions for u_i with an arbitrariness corresponding to a small deformation of the sheet in which the lengths of cord elements are unaltered.

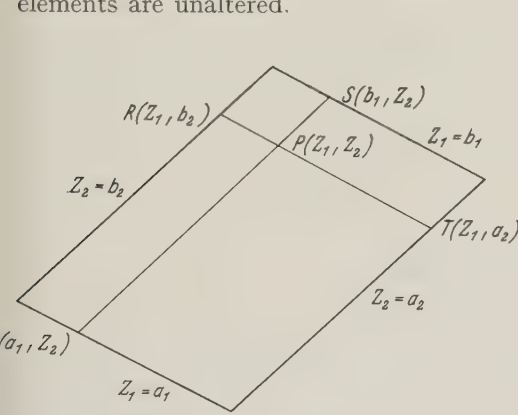


Fig. 3. Plane sheet bounded by four cords

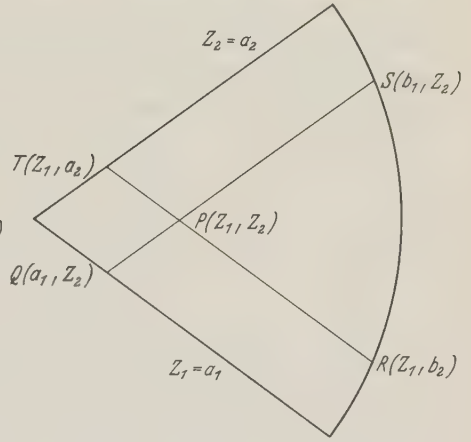


Fig. 4. Plane sheet bounded by two cords and an arc

(ii) *Sheet bounded by two intersecting cords and an arc.* We consider the boundary C of the sheet to be formed by the cords $Z_1 = a_1$ and $Z_2 = a_2$ and by an arc which is nowhere parallel to a cord, as shown in Fig. 4. Again P is a generic point of the sheet and Q, R, S, T are the intersections of the cords through P with the boundary. Suppose that S and R are the points (b_1, Z_2) and (Z_1, b_2) respectively.

From (9.3) and (9.4), we have

$$\begin{aligned}
 e_2(Z_1) &= \frac{\bar{d}}{k \cos \alpha} \xi_1(Z_1, a_2) = -\frac{\bar{d}}{k \sin \alpha} \xi_2(Z_1, a_2) \quad \text{at } T, \\
 e_1(Z_2) &= -\frac{\bar{d}}{k \cos \alpha} \xi_1(a_1, Z_2) = -\frac{\bar{d}}{k \sin \alpha} \xi_2(a_1, Z_2) \quad \text{at } Q.
 \end{aligned} \tag{9.13}$$

From (9.1), we have

$$\begin{aligned}
 e_1(Z_2) &= \frac{\bar{d} [\xi_1(b_1, Z_2) \sin \alpha + \xi_2(b_1, Z_2) \cos \alpha]}{k [M_2(b_1, Z_2) \cos \alpha - M_1(b_1, Z_2) \sin \alpha]} \quad \text{at } S \\
 e_2(Z_1) &= \frac{\bar{d} [\xi_1(Z_1, b_2) \sin \alpha - \xi_2(Z_1, b_2) \cos \alpha]}{k [M_2(Z_1, b_2) \cos \alpha + M_1(Z_1, b_2) \sin \alpha]} \quad \text{at } R.
 \end{aligned} \tag{9.14}$$

Employing either (9.13) or (9.14) in (7.8), we can obtain the displacement u_i throughout the sheet, with an arbitrariness corresponding to any deformation in which lengths of cord elements are unchanged.

We note that if the edge tractions are specified on the bounding cords, then it is sufficient that only one of the components ξ_i be specified on each of these cords. On the other hand, if the edge traction is specified on the bounding arc then both of the components ξ_i must be specified.

10. Mixed Boundary Value Problems

In this section we shall discuss a few examples which illustrate the manner in which the results obtained in the previous sections can be used to determine the displacements throughout a sheet when the boundary conditions are of the mixed type, *i.e.* displacements are given on parts of the boundary and edge tractions on other parts.

(i) *Sheet bounded by cords.* We consider the sheet of material bounded by four cords as shown in Fig. 3, and assume that the displacement component u_1 is specified on the edges $Z_1 = b_1$ and $Z_2 = b_2$, while the component of edge traction ξ_1 is specified on $Z_1 = a_1$ and $Z_2 = a_2$. From (9.3) and (9.4), we have

$$\begin{aligned} e_2(Z_1) &= \frac{\bar{d}}{k \cos \alpha} \xi_1(Z_1, a_2) \quad \text{at } T \\ \text{and} \\ e_1(Z_2) &= -\frac{\bar{d}}{k \cos \alpha} \xi_1(a_1, Z_2) \quad \text{at } Q. \end{aligned} \quad (10.1)$$

From the first of equations (7.8), we obtain

$$\begin{aligned} 2 \cos \alpha u_1(b_1, Z_2) &= b_1 e_1(Z_2) + Z_2 e_2(b_1) + p_1(Z_2) + p_2(b_1) \quad \text{at } S, \\ 2 \cos \alpha u_1(Z_1, b_2) &= Z_1 e_1(b_2) + b_2 e_2(Z_1) + p_1(b_2) + p_2(Z_1) \quad \text{at } R \\ \text{and} \\ 2 \cos \alpha u_1(b_1, b_2) &= b_1 e_1(b_2) + b_2 e_2(b_1) + p_1(b_2) + p_2(b_1) \quad \text{at } Z_1 = b_1, Z_2 = b_2. \end{aligned} \quad (10.2)$$

Employing the first two of equations (10.2) to eliminate $p_1(Z_2)$ and $p_2(Z_1)$ from (7.8), we obtain

$$\begin{aligned} 2 \cos \alpha u_1(Z_1, Z_2) &= 2 \cos \alpha [u_1(Z_1, b_2) + u_1(b_1, Z_2)] + \\ &\quad + (Z_1 - b_1) e_1(Z_2) + (Z_2 - b_2) e_2(Z_1) - \\ &\quad - Z_1 e_1(b_2) - Z_2 e_2(b_1) - p_1(b_2) - p_2(b_1) \\ \text{and} \\ 2 \sin \alpha u_2(Z_1, Z_2) &= 2 \cos \alpha [u_1(b_1, Z_2) - u_1(Z_1, b_2)] + \\ &\quad + (Z_1 - b_1) e_1(Z_2) - (Z_2 - b_2) e_2(Z_1) + \\ &\quad + Z_1 e_1(b_2) - Z_2 e_2(b_1) + p_1(b_2) - p_2(b_1). \end{aligned} \quad (10.3)$$

Eliminating $p_1(b_2) + p_2(b_1)$ from (10.3) and the third of equations (10.2), we obtain

$$\begin{aligned} 2 \cos \alpha u_1(Z_1, Z_2) &= 2 \cos \alpha [u_1(b_1, Z_2) + u_1(Z_1, b_2) - u_1(b_1, b_2)] + \\ &\quad + (Z_1 - b_1) [e_1(Z_2) - e_1(b_2)] + \\ &\quad + (Z_2 - b_2) [e_2(Z_1) - e_2(b_1)]. \end{aligned} \quad (10.4)$$

From (10.1), we have

$$\begin{aligned} e_2(b_1) &= \frac{\bar{d}}{k \cos \alpha} \xi_1(b_1, a_2) \\ \text{and} \\ e_1(b_2) &= -\frac{\bar{d}}{k \cos \alpha} \xi_1(a_1, b_2). \end{aligned} \quad (10.5)$$

Introducing (10.4) and (10.5) into (10.4) and the second of equations (10.3), we obtain

$$u_1(Z_1, Z_2) = u_1(b_1, Z_2) + u_1(Z_1, b_2) - u_1(b_1, b_2) - \frac{\bar{d}}{2k \cos^2 \alpha} \{ (Z_1 - b_1) [\xi_1(a_1, Z_2) - \xi_1(a_1, b_2)] - (Z_2 - b_2) [\xi_1(Z_1, a_2) - \xi_1(b_1, a_2)] \}$$

and

$$u_2(Z_1, Z_2) = \cot \alpha [u_1(b_1, Z_2) - u_1(Z_1, b_2)] - \frac{\bar{d}}{k \sin 2\alpha} [(Z_1 - b_1) \xi_1(a_1, Z_2) + (Z_2 - b_2) \xi_1(Z_1, a_2) + Z_1 \xi_1(a_1, b_2) + Z_2 \xi_1(b_1, a_2)] + \frac{\bar{p}_1(b_2) - \bar{p}_2(b_1)}{2 \sin \alpha}.$$

Thus, if u_1 is specified on $Z_1 = b_1$ and $Z_2 = b_2$ and ξ_1 is specified on $Z_1 = a_1$ and $Z_2 = a_2$, then u_1 is determined throughout the sheet and u_2 is determined throughout the sheet apart from a constant. If, further, u_2 is specified at a single point of the sheet, then u_2 is completely determined throughout the sheet.

(ii) *Sheet bounded by two cords and an arc. Prescribed tractions on the arc and displacements on the cords.* We consider a sheet of material bounded by the cords $Z_1 = a_1$ and $Z_2 = a_2$ and by an arc which is nowhere parallel to a cord, as shown in Fig. 4. We assume that u_1 is specified on the cords $Z_1 = a_1$ and $Z_2 = a_2$ and that the edge traction ξ_i is specified on the arc.

Then, from the first of equations (7.8), we obtain

$$\begin{aligned} 2 \cos \alpha u_1(Z_1, a_2) &= Z_1 e_1(a_2) + a_2 e_2(Z_1) + \bar{p}_1(a_2) + \bar{p}_2(Z_1) \quad \text{at } T, \\ 2 \cos \alpha u_1(a_1, Z_2) &= a_1 e_1(Z_2) + Z_2 e_2(a_1) + \bar{p}_1(Z_2) + \bar{p}_2(a_1) \quad \text{at } Q \\ 2 \cos \alpha u_1(a_1, a_2) &= a_1 e_1(a_2) + a_2 e_2(a_1) + \bar{p}_1(a_2) + \bar{p}_2(a_1) \quad \text{at } Z_1 = a_1, Z_2 = a_2. \end{aligned}$$

From (10.7) and (7.8), we obtain, in a manner similar to that used in discussing case (i),

$$\begin{aligned} 2 \cos \alpha u_1(Z_1, Z_2) &= 2 \cos \alpha [u_1(a_1, Z_2) + u_1(Z_1, a_2) - u_1(a_1, a_2)] + \\ &\quad + (Z_1 - a_1) [e_1(Z_2) - e_1(a_2)] + \\ &\quad + (Z_2 - a_2) [e_2(Z_1) - e_2(a_1)] \end{aligned}$$

and

$$\begin{aligned} 2 \sin \alpha u_2(Z_1, Z_2) &= 2 \cos \alpha [u_1(a_1, Z_2) - u_1(Z_1, a_2)] + \\ &\quad + (Z_1 - a_1) e_1(Z_2) - (Z_2 - a_2) e_2(Z_1) + \\ &\quad + Z_1 e_1(a_2) - Z_2 e_2(a_1) + \bar{p}_1(a_2) - \bar{p}_2(a_1). \end{aligned}$$

Now, solving equations (7.10) for $e_1(Z_2)$ and $e_2(Z_1)$ and employing (4.7), we obtain

$$\begin{aligned} e_1(Z_2) &= \frac{\bar{d} [\xi_1(Z_1, Z_2) \sin \alpha + \xi_2(Z_1, Z_2) \cos \alpha]}{k [M_2(Z_1, Z_2) \cos \alpha - M_1(Z_1, Z_2) \sin \alpha]} \\ e_2(Z_1) &= \frac{\bar{d} [\xi_1(Z_1, Z_2) \sin \alpha - \xi_2(Z_1, Z_2) \cos \alpha]}{k [M_2(Z_1, Z_2) \cos \alpha + M_1(Z_1, Z_2) \sin \alpha]}. \end{aligned}$$

In (10.9), the denominators are not zero on the bounding arc, since the arc is nowhere parallel to a cord. From (10.9) we have

$$\begin{aligned} e_1(Z_2) &= \frac{\bar{d} [\xi_1(b_1, Z_2) \sin \alpha + \xi_2(b_1, Z_2) \cos \alpha]}{k [M_2(b_1, Z_2) \cos \alpha - M_1(b_1, Z_2) \sin \alpha]} \quad \text{at } S \\ e_2(Z_1) &= \frac{\bar{d} [\xi_1(Z_1, b_2) \sin \alpha - \xi_2(Z_1, b_2) \cos \alpha]}{k [M_2(Z_1, b_2) \cos \alpha + M_1(Z_1, b_2) \sin \alpha]} \quad \text{at } R. \end{aligned}$$

If the points of intersection of $Z_1=a_1$ and $Z_2=a_2$ with the arc are denoted (a_1, c_2) and (c_1, a_2) respectively, then from (10.9) we obtain

$$\begin{aligned} e_1(a_2) &= \frac{\bar{d}[\xi_1(c_1, a_2) \sin \alpha + \xi_2(c_1, a_2) \cos \alpha]}{k[M_2(c_1, a_2) \cos \alpha - M_1(c_1, a_2) \sin \alpha]} \\ e_2(a_1) &= \frac{\bar{d}[\xi_1(a_1, c_2) \sin \alpha - \xi_2(a_1, c_2) \cos \alpha]}{k[M_2(a_1, c_2) \cos \alpha + M_1(a_1, c_2) \sin \alpha]}. \end{aligned} \quad (10.11)$$

Employing equations (10.10) and (10.11) to substitute for $e_1(Z_2)$, $e_2(Z_1)$, $e_1(a_2)$ and $e_2(a_1)$ in (10.8), we obtain expressions for $u_i(Z_1, Z_2)$ in terms of the known value of u_1 on the bounding cords $Z_1=a_1$ and $Z_2=a_2$ and the known values of the edge tractions ξ_i on the bounding arc. We note that u_2 is determined to within an arbitrary constant. However, if u_2 is prescribed at a single point of the sheet, it is then completely determined throughout the sheet.

(iii) *Sheet bounded by two cords and an arc. Prescribed tractions on the cords and displacements on the arc.* We again consider the sheet shown in Fig. 4, but we now assume that the displacement components u_i are known on the arc and the component of edge traction ξ_1 is known on the cords $Z_1=a_1$ and $Z_2=a_2$.

From (9.3) and (9.4) we obtain

$$\begin{aligned} e_2(Z_1) &= -\frac{\bar{d}}{k \cos \alpha} \xi_1(Z_1, a_2) \quad \text{at } T \\ e_1(Z_2) &= -\frac{\bar{d}}{k \cos \alpha} \xi_1(a_1, Z_2) \quad \text{at } Q. \end{aligned} \quad (10.12)$$

From (7.7) it is seen that

$$\begin{aligned} \cos \alpha u_1(Z_1, b_2) - \sin \alpha u_2(Z_1, b_2) &= b_2 e_2(Z_1) + p_2(Z_1) \quad \text{at } R \\ \cos \alpha u_1(b_1, Z_2) + \sin \alpha u_2(b_1, Z_2) &= b_1 e_1(Z_2) + p_1(Z_2) \quad \text{at } S. \end{aligned} \quad (10.13)$$

Employing (10.13) to eliminate $p_1(Z_2)$ and $p_2(Z_1)$ from (7.8), we have

$$\begin{aligned} 2 \cos \alpha u_1(Z_1, Z_2) &= \cos \alpha [u_1(b_1, Z_2) + u_1(Z_1, b_2)] + \\ &\quad + \sin \alpha [u_2(b_1, Z_2) - u_2(Z_1, b_2)] + \\ &\quad + (Z_1 - b_1) e_1(Z_2) + (Z_2 - b_2) e_2(Z_1) \\ 2 \sin \alpha u_2(Z_1, Z_2) &= \cos \alpha [u_1(b_1, Z_2) - u_1(Z_1, b_2)] + \\ &\quad + \sin \alpha [u_2(b_1, Z_2) + u_2(Z_1, b_2)] + \\ &\quad + (Z_1 - b_1) e_1(Z_2) - (Z_2 - b_2) e_2(Z_1). \end{aligned} \quad (10.14)$$

Substituting for $e_2(Z_1)$ and $e_1(Z_2)$ from (10.12) in (10.14) we obtain expressions for $u_i(Z_1, Z_2)$ in terms of the known component of edge traction ξ_1 on $Z_1=a_1$ and $Z_2=a_2$ and the known displacement components on the bounding arc.

Acknowledgement. The results presented in this paper were obtained in the course of research sponsored by the Office of Ordnance Research, U. S. Army, under Contract No. DA-19-020-ORD-4725 and its precursors.

Brown University
Providence, Rhode Island

(Received May 15, 1959)

Iterationsverfahren mit veränderlichen Operatoren

HANS EH RMANN

Vorgelegt von L. COLLATZ

§ 1. Einleitung^{*}

Iterationsverfahren der Gestalt

$$(1.1) \quad u_{n+1} = T u_n, \quad n = 0, 1, 2, \dots,$$

zur Lösung einer Gleichung

$$(1.2) \quad u = T u \quad \text{oder} \quad S u = \Theta, \quad (\Theta = \text{Nullelement}),$$

wobei die Größen $u, u_n, n=0, 1, 2, \dots$, Elemente eines abstrakten Raumes (z.B. eines Banachraumes) sind und T einen Operator bedeutet, sind in den letzten Jahren viel in bezug auf Konvergenzfragen, Konvergenzverbesserungen, Fehlerabschätzungen usw. untersucht worden^{**}. Dabei zeigte es sich, daß viele bei speziellen Gleichungstypen bekannte Konvergenzkriterien und Fehlerabschätzungen, die vielfach in jedem speziellen Falle neu bewiesen wurden, Spezialfälle von z.T. sehr einfachen und leicht zu beweisenden abstrakten Sätzen sind und daß sich die Verschiedenheit von speziellen Konvergenzbedingungen und Fehlerabschätzungen aus der verschiedenen speziellen Wahl der Metrik ergibt [3], [5], [10]. In den meisten Fällen beruht die Konvergenz bei den betr. Sätzen auf der Eigenschaft des Operators T zu kontrahieren, und im einzelnen wird gewöhnlich gezeigt, daß die Summe der Abstände von u_{n+1} und u_n für $n=0, 1, 2, \dots$ konvergiert, was i. allg. mit Hilfe des Majorantenprinzips mittels einer Vergleichsreihe mit positiven Elementen geschieht, deren Konvergenz bekannt ist. Durch die abstrakte Formulierung der Sätze gelingt es zunächst,

^{*} Dies ist ein erster Auszug einer Arbeit „Untersuchungen von Iterationsverfahren zur Lösung allgemeiner Gleichungen“, die im November 1958 von der Fakultät für Natur- und Geisteswissenschaften der Technischen Hochschule Stuttgart als Habilitationsschrift anerkannt wurde. Die Untersuchungen entstanden während meiner Tätigkeit als Assistent am Institut für Mathematik und Mechanik der Bergakademie Clausthal bei Herrn Professor H. KÖNIG und am Mathematischen Institut A der Technischen Hochschule Stuttgart bei Herrn Professor G. SCHULZ. Daneben hat Herr Professor L. COLLATZ die Arbeit verfolgt und mich durch einige wertvolle Ratschläge und Hinweise unterstützt. Allen drei Professoren möchte ich für ihre freundliche Unterstützung herzlich danken.

^{**} Siehe u. a. das Literaturverzeichnis am Ende der Arbeit. Hierauf beziehen sich die Literaturangaben in eckigen Klammern.

einen wesentlich größeren Kreis von Gleichungstypen (Gleichungssysteme, Differential- und Integralgleichungen usw.) zu erfassen. In der weiteren Entwicklung wurde sodann versucht, Konvergenzbedingungen und Fehlerabschätzungen der betr. abstrakten Sätze zu verbessern, indem die Definitionen insbes. der Metrik weiter verallgemeinert wurden, um die Abschätzungen dem jeweiligen speziellen Problem besser anpassen zu können. Dabei werden auch metrische Räume mit „Abständen“ ihrer Punkte zugelassen, die keine nichtnegativen reellen Zahlen sind, [2], [7], [9]. Solche allgemeineren Räume, die den Banachraum als Spezialfall enthalten, werden wir auch hier zugrunde legen*.

Es treten nun Fälle auf, in denen z.B. die Iteration (1.1) noch sehr mühsam ist, so daß man zweckmäßig wenigstens in den ersten Schritten T durch einen Näherungsoperator \tilde{T} ersetzt, oder T ist von vornherein nur näherungsweise gegeben etwa in Form einer unendlichen Reihe, von der in (1.1) jeweils nur endlich viele Glieder berücksichtigt werden können, oder die Lösung u erscheint in Form einer unendlichen Reihe wie z.B. bei dem bekannten Verfahren zur Auflösung einer Gleichung der Form $x=f(x, \lambda)$, wobei man f an einer Stelle x_0 nach Potenzen von λ entwickelt, die Reihe abbricht, rechts für x einsetzt, wieder entwickelt, abbricht, einsetzt usw.:

$$x_{n+1} = f(x_n, \lambda) - \sum_{v \geq k_n} a_{v,n} \lambda^v \quad \text{mit} \quad f(x_n, \lambda) = \sum_{v=0}^{\infty} a_{v,n} \lambda^v.$$

Diese Fragestellungen führen auf das Problem, Iterationsverfahren der Form

$$(1.3) \quad u_{n+1} = T_n u_n$$

zu untersuchen, wobei sich also die Operatoren T_n von Schritt zu Schritt ändern können. Hiermit beschäftigt sich die vorliegende Arbeit. Es wird ein abstrakter Satz aufgestellt und bewiesen, der Bedingungen angibt, unter denen die Iteration (1.3) gegen eine Lösung der Gleichung (1.1) konvergiert. Außerdem wird eine Fehlerabschätzung angegeben, die unter den angegebenen Bedingungen nicht mehr verbessert werden kann. Der entsprechende Satz geht in dem Spezialfall $T_n = T$, $n=1, 2, \dots$, in bekannte Sätze über [4], [9]. Jedoch versagt hier die übliche Beweismethode. So braucht z.B. unter den getroffenen Voraussetzungen die Summe der Abstände von den aufeinanderfolgenden Iterationspunkten u_{n+1} und u_n ($n=0, 1, 2, \dots$) nicht zu konvergieren, noch kann man zeigen, daß diese Abstände schließlich monoton abnehmen. Obwohl der Satz dem Inhalt nach ganz elementar ist, ist er meines Wissens auch im einfachsten Fall einer unbekannten Zahl noch nicht in dieser Form mit Konvergenzbeweis und Fehlerabschätzung bekannt. Andererseits ist die hier durchgeführte abstrakte Formulierung so allgemein, daß der Satz auf sehr verschiedenartige Gleichungen leicht angewandt werden kann.

In § 2 werden alle funktionalanalytischen Hilfsmittel zusammengestellt. § 3 bringt den Satz und Beweis sowie einige für die Anwendungen auf den Nachweis der Existenz von Lösungen und ihre Abschätzungen wichtige Folgerungen. In § 4 werden schließlich einige Beispiele durchgeführt.

* Vgl. insbes. J. SCHRÖDER [9].

§ 2. Funktionalanalytische Voraussetzungen

Für je zwei Elemente $u, v \in R$ sei ein Abstand $\varrho(u, v)$ definiert, wobei die Abstände ϱ, σ, \dots Elemente eines linearen halbgeordneten Raumes N sind: es ist in N eine Addition $\varrho + \sigma$ und eine Multiplikation $\alpha \varrho$ mit reellen Zahlen α definiert, und diese Operationen genügen den üblichen Rechenregeln der Vektoralgebra. Ferner ist N halbgeordnet, d.h. es sind gewisse Elemente $\varrho, \sigma \in N$ als positiv definiert: $\varrho \geq 0$ ($0 =$ Nullelement von N), wobei $\varrho \geq 0$ und $-\varrho \geq 0$ (gleichzeitig) nur für $\varrho = 0$ gilt und $\varrho + \sigma \geq 0$ aus $\varrho, \sigma \geq 0$, sowie $\alpha \varrho \geq 0$ aus $\alpha \geq 0$ (α reelle Zahl) und $\varrho \geq 0$ folgt. $\varrho \geq \sigma$ oder $\sigma \leq \varrho$ bedeutet $\varrho - \sigma \geq 0$.

Ferner ist gewissen (konvergent genannten) Folgen $\varrho_n \in N$ ($n = 1, 2, 3, \dots$) ein (eindeutig bestimmtes) Grenzelement $\lim \varrho_n \in N$ zugeordnet mit den Eigenschaften*:

- a) aus $\varrho_n = \varrho$ für $n = 1, 2, 3, \dots$ mit einem festen Element ϱ folgt $\lim \varrho_n = \varrho$,
 b) aus $\lim \varrho_n = \varrho$ folgt $\lim \varrho_{k_n} = \varrho$ für jede monotone Folge natürlicher Zahlen k_n ($n = 1, 2, 3, \dots$),
 c) aus $\lim \varrho_n = \varrho$ und $\lim \sigma_n = \sigma$ folgt $\lim (\varrho_n + \sigma_n) = \varrho + \sigma$,
 (2.1) d) aus $\lim \alpha_n = \alpha$ (α_n, α reelle Zahlen) und $\lim \varrho_n = \varrho$ folgt $\lim \alpha_n \varrho_n = \alpha \varrho$,
 e) aus $0 \leq \varrho_n \leq \sigma_n$ für $n = 1, 2, 3, \dots$ und $\lim \sigma_n = 0$ folgt $\lim \varrho_n = 0$,
 f) aus $\varrho_n \geq 0$ für $n = 1, 2, 3, \dots$ und $\lim \varrho_n = \varrho$ folgt $\varrho \geq 0$,
 g) aus $0 \leq \varrho_n \leq \tau$ für $n \geq n_0$ und $\lim \varrho_n = 0$ folgt die Existenz einer Nummer n_1 mit $\varrho_n \leq \frac{\tau}{2}$ für $n \geq n_1$.

Ferner gelten für den bezüglich N metrischen Raum R die Abstandspostulate:

1. $\varrho(u, v) = 0$ genau dann, wenn $u = v$ gilt,
2. $\varrho(u, v) \leq \varrho(u, w) + \varrho(v, w)$ für jedes Elemententripel $u, v, w \in R$ (Dreiecksungleichung).

Daraus folgt dann die Symmetrie, $\varrho(u, v) = \varrho(v, u)$, und die Definitheit, $\varrho(u, v) \geq 0$, des Abstandes. Man beweist leicht, daß die Axiome (2.1e) und (2.1g) in dem obigen System äquivalent sind den Axiomen:

(2.1e') Aus $\tau_n \leq \varrho_n \leq \sigma_n$ für $n = 1, 2, \dots$ und $\lim \tau_n = \lim \sigma_n = \sigma$ folgt $\lim \varrho_n = \sigma$ bzw.

(2.1g') aus $0 \leq \varrho_n \leq \tau$ für $n \geq n_0$ und $\lim \varrho_n = 0$ folgt für jede natürliche Zahl k die Existenz einer Nummer $N(k)$ mit $\varrho_n \leq \frac{\tau}{k}$ für $n \geq N(k)$.

Die Konvergenz einer Folge $u_n \rightarrow u$, $n = 0, 1, 2, \dots$ in R ist stets im Sinne der sog. starken Konvergenz bezüglich der Metrik N zu verstehen, d.h.

$$(2.2) \quad u = \lim_{n \rightarrow \infty} u_n \text{ ist äquivalent } \lim_{n \rightarrow \infty} \varrho(u_n, u) = 0.$$

* Die Axiome a) bis f) sind wörtlich von J. SCHRÖDER [9] übernommen. Zum Beweis des in § 3 aufgestellten Satzes über Iterationsverfahren mit veränderlichen Operatoren (1.3) erwies sich das zusätzliche Axiom g) als erforderlich, dessen Unabhängigkeit von den übrigen etwa durch das Beispiel $\varrho_n = f_n(x) = \frac{nx}{1+n^2x^2}$ und $\tau = \frac{1}{2}$ im Raum der stetigen Funktionen gezeigt wird.

Man zeigt leicht:

$$(2.3) \quad \text{Aus } u_n \rightarrow u, \quad v_n \rightarrow v \quad \text{folgt} \quad \varrho(u_n, v_n) \rightarrow \varrho(u, v).$$

Außerdem wird R als *vollständig* vorausgesetzt, d.h. für jede Folge $\{u_n\}$ in R , für die das Cauchysche Konvergenzkriterium $\lim_{l,k \rightarrow \infty} \varrho(u_l, u_k) = 0$ gilt, existiert ein Grenzelement $u = \lim_{n \rightarrow \infty} u_n$ in R .

Der bezüglich N metrische Raum R heißt *lokal kompakt*, wenn jede unendliche Folge $u_n \in R$, $n = 0, 1, 2, \dots$, für die $\varrho(u_n, v) \leq \sigma$, $n = 0, 1, 2, \dots$, mit $v \in R$ und $\sigma \in N$ gilt, eine konvergente Teilfolge u_{k_n} besitzt.

Eine Folge von Operatoren T_n , $n = 0, 1, 2, \dots$, die in einem Gebiet D von R definiert sind, konvergiert dort *stetig*^{*} gegen einen Operator T , wenn für jede Folge $\{v_n\}$ aus D , die gegen einen Punkt v von D konvergiert, der Grenzwert $\lim_{n \rightarrow \infty} T_n v_n$ existiert und gleich Tv ist.

Eine Folge von Operatoren T_n , $n = 0, 1, 2, \dots$, konvergiert in einem Bereich G ihres Definitionsbereiches *gleichmäßig* gegen T , wenn $\varrho(T_n u, T u) \leq \sigma$ für alle $u \in G$ mit einer Größe $\sigma \in N$ ist und nach Vorgabe einer positiven Zahl k eine Nummer $n(k)$ existiert derart, daß

$$\varrho(T_n u, T u) \leq \frac{\sigma}{k} \quad \text{für } n \geq n(k) \text{ und alle } u \in G$$

gilt.

Ein in einem Bereich G eines bezüglich N metrischen Raumes R definierter Operator T heißt dort *beschränkt*, wenn

$$(2.4) \quad \varrho(Tu, Tv) \leq P\varrho(u, v), \quad u, v \in G,$$

mit einem linearen, stetigen und positiven Operator P gilt. Dabei heißt P *linear*, wenn $P(\alpha\varrho + \beta\sigma) = \alpha P\varrho + \beta P\sigma$ (α, β reelle Zahlen, $\varrho, \sigma \in N$) gilt, *stetig*, wenn $\lim P\varrho_n = P\varrho$ aus $\lim \varrho_n = \varrho$, und *positiv*, wenn $P\varrho \geq 0$ aus $\varrho \geq 0$ folgt.

Ein linearer positiver Operator P ist auch monoton, d.h. es gilt

$$(2.5) \quad \text{aus } \varrho \geq \sigma, \quad P \text{ linear, positiv, folgt } P\varrho \geq P\sigma.$$

Schließlich verwenden wir noch folgenden einfachen

Hilfssatz 1: Für die Folge von Operatoren T_n , $n = 0, 1, 2, \dots$ gelte in einem Bereich G ihres gemeinsamen Definitionsbereiches $D \subseteq R$

$$(2.6) \quad \lim_{n \rightarrow \infty} \varrho(Tu, T_n u) = 0, \quad u \in G,$$

und

$$(2.7) \quad \varrho(T_n u, T_n v) \leq P\varrho(u, v), \quad n = 0, 1, 2, \dots, \quad u, v \in G,$$

mit einem stetigen Operator P . Dann konvergieren die Operatoren T_n , $n = 0, 1, 2, \dots$, stetig in G .

Beweis. Es gelte $\lim_{n \rightarrow \infty} v_n = v$, $v_n, v \in G$. Dann folgt nach der Dreiecksungleichung und (2.7):

$$0 \leq \varrho(T_n v_n, Tv) \leq \varrho(T_n v_n, T_n v) + \varrho(T_n v, Tv) \leq P\varrho(v_n, v) + \varrho(T_n v, Tv).$$

Hieraus ergibt sich wegen (2.6) und der Stetigkeit von P mit (2.1c) und (2.1e) die Behauptung.

* Vgl. C. CARATHÉODORY [13] (S. 1ff.).

§ 3. Der Fixpunktsatz

Wir treffen folgende *allgemeine Voraussetzungen*:

A. Es sei R ein bezüglich eines linearen halbgeordneten Raumes N metrischer, vollständiger Raum^{*}. Dabei sei die Konvergenz in N durch die Axiome (2.1 a—g) festgelegt.

B. Ferner sei eine Folge von in einem gemeinsamen Teilgebiet D von R definierten Operatoren

$$(3.1) \quad T_0, T_1, T_2, \dots$$

gegeben mit den Eigenschaften:

1 a. Die Folge (3.1) konvergiert gleichmäßig^{*, **} in D gegen einen Operator T

$$(3.2) \quad \lim_{n \rightarrow \infty} \varrho(T_n u, T u) = 0, \quad u \in D.$$

Statt der gleichmäßigen Konvergenz können wir fordern:

1 a'. Der Raum R ist lokalkompakt^{*} und es gilt (3.2) ohne die Bedingung (2.1 g).

2. Die Operatoren (3.1) seien in einem beschränkten Teilgebiet G von D gleichmäßig beschränkt, d.h. es existiert ein stetiger, positiver, linearer Operator P , so daß

$$(3.3) \quad \varrho(T_n u, T_n v) \leq P \varrho(u, v), \quad n = 0, 1, 2, \dots, \quad u, v \in G \subseteq D,$$

gilt.

3. Schließlich existiere eine obere Schranke $\mu \in N$ derart, daß an einer Stelle u_0 von G

$$(3.4) \quad \varrho(T_0 u_0, T_n u_0) \leq \mu, \quad n = 0, 1, 2, \dots, \quad u_0 \in G,$$

gilt.

Es ergeben sich einige einfache Folgerungen aus A und B:

Zunächst folgt aus (3.2) und (3.3), daß auch der Grenzoperator T beschränkt ist mit demselben Operator P , d.h. es gilt

$$(3.5) \quad \varrho(T u, T v) \leq P \varrho(u, v), \quad u, v \in G.$$

Dies folgt aus (2.3) und den Axiomen (2.1 a, c, d, f) auf die Folge

$$\tau_n = P \varrho(u, v) - \varrho(T_n u, T_n v) \geq 0, \quad n = 0, 1, 2, \dots,$$

angewandt.

Aus (3.3) und (3.4) folgt für das gesamte beschränkte Gebiet G die Existenz einer oberen Schranke für $\varrho(T_n u, T_n u)$: Zunächst folgt aus der Dreiecksungleichung (3.3) und (3.4) für $u \in G$

$$\begin{aligned} \varrho(T_0 u, T_n u) &\leq \varrho(T_0 u, T_0 u_0) + \varrho(T_0 u_0, T_n u_0) + \varrho(T_n u_0, T_n u) \\ &\leq 2P \varrho(u, u_0) + \mu \leq \tilde{\mu}, \quad n = 0, 1, 2, \dots, \end{aligned}$$

^{*} Bezügl. der Definitionen siehe § 2.

^{**} Später werden wir uns noch von dieser Voraussetzung befreien.

denn da G beschränkt ist, liegt $\varrho(u, u_0)$ und damit auch $2P\varrho(u, u_0)$ unter einer festen Schranke. Weiter folgt aus der letzten Gleichung

$$(3.6) \quad \varrho(T_m u, T_n u) \leq \varrho(T_0 u, T_m u) + \varrho(T_0 u, T_n u) \leq 2\tilde{\mu}, \\ n, m = 0, 1, 2, \dots, \quad u \in G.$$

Läßt man hierin $n \rightarrow \infty$ gehen, so ergibt sich nach einfacher Überlegung wegen der Stetigkeit des Abstandes [(2.3)] auch

$$(3.7) \quad \varrho(T_m u, T u) \leq 2\tilde{\mu}, \quad m = 0, 1, 2, \dots, \quad u \in G.$$

Schließlich folgt nach Hilfssatz 1 noch, daß die Folge der Operatoren T_n in G stetig gegen T konvergiert.

Unter den allgemeinen Voraussetzungen A und B gilt der

Fixpunktsatz: Es gelte*

$$\tau_0 \geq \varrho(T_n u_0, T_0 u_0), \quad n = 0, 1, 2, \dots$$

Liegen u_0 und alle u , für die**

$$(3.8) \quad \varrho(u, u_1) \leq (E - P)^{-1} P \varrho(u_0, u_1) + (E - P)^{-1} \tau_0 = \sigma$$

gilt, in dem Gültigkeitsbereich G der Gleichung (3.3) und konvergiert die Neumannsche Reihe

$$(3.9) \quad \sum_{\nu=0}^{\infty} P^{\nu} \varrho \quad \text{für alle } \varrho \in N,$$

so konvergiert das Iterationsverfahren

$$(3.10) \quad u_{n+1} = T_n u_n$$

gegen die einzige Lösung u der Gleichung

$$u = T u$$

in G und es gilt die Fehlerabschätzung***

$$(3.11) \quad \varrho(u, u_1) \leq (E - P)^{-1} P \varrho(u_0, u_1) + (E - P)^{-1} \varrho(T u_0, u_1).$$

Beweis. $\alpha)$ Mit P ist** auch $(E - P)^{-1}$ ein positiver Operator. Daher ist die rechte Seite von (3.8) $\sigma \geq 0$. Es liegt also u_1 in dem durch (3.8) gegebenen Bereich von G . Wir treffen die Induktionsannahme, daß u_1, u_2, \dots, u_s in G liegen. Dann sei für $n=1, 2, \dots, s+1$

$$(3.12) \quad \tau_n \geq \varrho(T_l u_n, T_n u_n) \quad \text{für } l \geq n.$$

Die Existenz dieser oberen Schranken τ_n folgt aus (3.4) und (3.6).

* Die Existenz einer solchen oberen Schranke ergibt (3.4), z. B. $\tau_0 = \mu$. Existiert eine kleinste obere Schranke, so wird man zweckmäßig $\tau_0 = \sup_n \varrho(T_n u_0, T_0 u_0)$ setzen.

** Die Existenz des Operators $(E - P)^{-1}$ folgt aus der Konvergenz der Reihe (3.9), und es gilt $(E - P)^{-1} \varrho = \sum_{\nu=0}^{\infty} P^{\nu} \varrho$, $\varrho \in N$, siehe J. SCHRÖDER [9].

*** Diese Fehlerabschätzung läßt sich natürlich auch mit u_{n+1} und u_n statt u_1 und u_0 schreiben, was nur eine andere Zählung bedeutet.

Sodann beweisen wir durch vollständige Induktion

$$(3.13) \quad \varrho(u_{n+k}, u_n) \leq (E - P)^{-1} P \varrho(u_n, u_{n-1}) + (E - P)^{-1} \tau_{n-1} = \sigma_n$$

für $1 \leq n \leq s$ und $k \geq 0$.

Es ist $\sigma_n \geq 0$ für alle n . Daher ist (3.13) für $k=0$ und $1 \leq n \leq s$ erfüllt. (3.13) sei für $k=0, 1, \dots, r$ erfüllt, und es mögen $u_{n-1}, u_n, u_{n+1}, \dots, u_{n+r}$ in dem Gültigkeitsbereich G von (3.3) liegen. Dann folgt aus (3.3), (3.12) und der Dreiecksungleichung

$$(3.14) \quad \begin{aligned} \varrho(u_{n+r+1}, u_n) &= \varrho(T_{n+r} u_{n+r}, T_{n-1} u_{n-1}) \\ &\leq \varrho(T_{n+r} u_{n+r}, T_{n+r} u_{n-1}) + \varrho(T_{n+r} u_{n-1}, T_{n-1} u_{n-1}) \\ &\leq P \varrho(u_{n+r}, u_{n-1}) + \tau_{n-1}. \end{aligned}$$

Da nach (2.5) der lineare, positive Operator P auch monoton ist, gilt nach der Dreiecksungleichung

$$P \varrho(u_{n+r}, u_{n-1}) \leq P \varrho(u_{n+r}, u_n) + P \varrho(u_n, u_{n-1}).$$

Setzt man dies in (3.14) ein und sodann für $\varrho(u_{n+r}, u_n)$ die rechte Seite von (3.13), was wegen der Induktionsvoraussetzung für k erlaubt ist, so ergibt sich

$$(3.15) \quad \varrho(u_{n+r+1}, u_n) \leq [P(E-P)^{-1} + E] P \varrho(u_n, u_{n-1}) + [P(E-P)^{-1} + E] \tau_{n-1}.$$

Nun ist aber

$$[P(E-P)^{-1} + E] \varrho = (E-P)^{-1} \varrho \quad \text{für alle } \varrho \in N,$$

denn es gilt

$$[P(E-P)^{-1} + E] \varrho - (E-P)^{-1} \varrho = [E - (E-P)(E-P)^{-1}] \varrho = [E-E] \varrho = 0.$$

Daher folgt aus (3.15) die Beziehung (3.13) für $k=r+1$, also unter den obigen Annahmen für alle $k \geq 0$ und $1 \leq n \leq s$. Setzt man hierin $n=1$, so erkennt man, da für u_0 und u_1 alle getroffenen Induktionsvoraussetzungen gültig sind, daß die Iteration (3.10) nicht aus dem durch (3.8) gegebenen Bereich herausführt. Es liegen damit alle u_n in G und daher folgt die Gültigkeit von (3.12) für alle $k \geq 0$ und $n \geq 1$.

β) Wir zeigen weiter, daß

$$(3.16) \quad \lim_{n \rightarrow \infty} \varrho(T_{l_n} u_{n-1}, T_{n-1} u_{n-1}) = 0$$

ist, wobei $\{l_n\}$ eine (nicht notwendig monotone) Folge ganzer Zahlen mit $l_n \geq n$ ist. Dies ist wegen

$$\varrho(T_{l_n} u_{n-1}, T_{n-1} u_{n-1}) \leq \varrho(T_{l_n} u_{n-1}, T u_{n-1}) + \varrho(T u_{n-1}, T_{n-1} u_{n-1})$$

der Fall, wenn die Größen

$$\varrho_m = \varrho(T_{k_m} u_m, T u_m) \quad \text{mit} \quad k_m \geq m$$

für $m \rightarrow \infty$ gegen Null konvergieren.

Dies ist aber offenbar wegen der gleichmäßigen Konvergenz (Voraussetzung B 1 a.) der Operatoren T_n , $n=0, 1, 2, \dots$, also auch der Operatoren T_{k_m} , $m=0, 1, 2, \dots$, $k_m \geq m$, gegen T der Fall, weil nach α) alle u_m in dem beschränkten Bereich (3.8) von G liegen.

Ist statt B 1 a. die Voraussetzung B 1 a'. erfüllt, so zeigen wir die Konvergenz der Folge $\varrho_m \rightarrow 0$ für $m \rightarrow \infty$ folgendermaßen:

Nach (3.7) und $u_m \in G$ für $m=0, 1, 2, \dots$ (siehe α)) liegen die Abstände ϱ_m alle unter einer festen Schranke τ . Würde ϱ_m nicht gegen Null konvergieren, so gäbe es eine positive Zahl p und eine Teilfolge ϱ_{r_n} , $n=1, 2, \dots$, für die die Gleichung $\varrho_{r_n} \leq \frac{\tau}{p}$ nicht erfüllt ist. Wir schreiben dafür *

$$(3.17) \quad \varrho_{r_n} \not\leq \frac{\tau}{p}, \quad n = 1, 2, \dots$$

Die unendliche Folge u_{r_n} besitzt wegen der Lokalkompaktheit des Teilraumes G von R eine konvergente Teilfolge. Wir können daher zugleich mit (3.17) annehmen, daß bereits u_{r_n} , $n=1, 2, \dots$, gegen ein Element v von G konvergiert. Wegen der stetigen Konvergenz der Operatoren T_n (Hilfssatz 1) konvergiert** damit jede Folge $T_n u_{r_n}$, $n=1, 2, \dots$, $l_n \geq r_n$, gegen Tv , und hieraus folgt mittels der Dreiecksungleichung

$$\lim_{n \rightarrow \infty} \varrho_{r_n} = 0.$$

Dieses steht aber nach dem Axiom (2.1 g') im Widerspruch zu (3.17). Hieraus folgt, daß die Größen ϱ_m gegen Null gehen mit $m \rightarrow \infty$, und damit ist auch die Behauptung (3.16) bewiesen.

γ) Wir wollen jetzt zeigen, daß die Folge $\{u_n\}$ gegen ein Element u konvergiert. Zu diesem Zweck beweisen wir zunächst die Ungleichung

$$(3.18) \quad \varrho(u_{n+r+k}, u_{n+k-1}) \leq P^k \varrho(u_{n+r}, u_{n-1}) + \sum_{v=1}^k P^{v-1} \tau_{n+k-v-1}, \quad n, k \geq 1,$$

wobei die τ_i als obere Schranken gemäß (3.12) definiert sind und $r=r(n, k)$ eine beliebige Doppelfolge natürlicher Zahlen bedeute.

Für $k=1$ erhalten wir nach (3.3), (3.12) und der Dreiecksungleichung

$$(3.19) \quad \begin{aligned} \varrho(u_{n+r+1}, u_n) &= \varrho(T_{n+r} u_{n+r}, T_{n-1} u_{n-1}) \\ &\leq \varrho(T_{n+r} u_{n+r}, T_{n+r} u_{n-1}) + \varrho(T_{n+r} u_{n-1}, T_{n-1} u_{n-1}) \\ &\leq P \varrho(u_{n+r}, u_{n-1}) + \tau_{n-1}, \quad n = 1, 2, \dots \end{aligned}$$

Also ist (3.18) für $k=1$ und alle n und r erfüllt. Es sei (3.18) für $k=1, 2, \dots, s$ erfüllt. Ersetzt man in (3.19) n durch $n+s$, so ergibt sich unter Benutzung von (3.18) für $k=s$

$$\begin{aligned} \varrho(u_{n+r+s+1}, u_{n+s}) &\leq P \varrho(u_{n+s+r}, u_{n+s-1}) + \tau_{n+s-1} \\ &\leq P^{s+1} \varrho(u_{n+r}, u_{n-1}) + P \sum_{v=1}^s P^{v-1} \tau_{n+s-v-1} + \tau_{n+s-1} \\ &\leq P^{s+1} \varrho(u_{n+r}, u_{n-1}) + \sum_{v=1}^{s+1} P^{v-1} \tau_{n+(s+1)-v-1}. \end{aligned}$$

Also ist (3.18) auch für $k=s+1$ und daher für alle $k \geq 1$ erfüllt.

Nach α) und (3.8) gilt

$$\varrho(u_m, u_n) \leq \varrho(u_m, u_1) + \varrho(u_n, u_1) \leq 2\sigma, \quad m, n = 1, 2, \dots$$

* Man beachte, daß aus $\varrho \not\leq \sigma$ noch nicht $\varrho \geq \sigma$ folgt.

** Mit T_n konvergiert auch die Folge T_{l_n} , $n=1, 2, \dots$, $l_n \geq n$, stetig, wie man leicht unter Verwendung des Axioms (2.1 b) schließt.

Daher ist auch $\varrho(u_{n+r}, u_{n-1}) \leq 2\sigma$ für $n \geq 2$, $r \geq 0$, und da P monoton und linear ist,

$$(3.20) \quad P^k \varrho(u_{n+r}, u_{n-1}) \leq 2P^k \sigma, \quad \sigma \in N.$$

Wegen der Konvergenz der Neumannschen Reihe für alle $\varrho \in N$ folgt, wie man leicht zeigen kann,

$$2P^k \sigma \rightarrow 0 \quad \text{mit} \quad k \rightarrow \infty,$$

und daher geht nach (2.1e) auch die linke Seite von (3.20), also der erste Summand der rechten Seite von (3.18), bei beliebigen $n \geq 1$ und $r \geq 0$ mit $k \rightarrow \infty$ gegen Null.

Nach β) konvergiert die rechte Seite von (3.12) gegen Null. Wir können daher wegen (2.1g') auch immer erreichen, daß die τ_n eine Nullfolge bilden*. Es seien nun t_n obere Schranken der τ_m für $m \geq n-1$. Dann können wir wegen $\tau_m \rightarrow 0$ für $m \rightarrow \infty$ wiederum nach (2.1g') annehmen, daß auch die $t_n \rightarrow 0$ gehen mit $n \rightarrow \infty$. Mit

$$t_n \geq \tau_{n+k-v-1}, \quad t_n \rightarrow 0 \quad \text{mit} \quad n \rightarrow \infty$$

ergibt sich für den zweiten Summanden der rechten Seite von (3.18)

$$\sum_{v=1}^k P^{v-1} \tau_{n+k-v-1} \leq \sum_{v=1}^k P^{v-1} t_n \leq \sum_{v=1}^{\infty} P^{v-1} t_n = (E - P)^{-1} t_n \rightarrow 0 \quad \text{mit} \quad n \rightarrow \infty.$$

Nach (2.1e) geht daher der zweite Summand von (3.18) mit $n \rightarrow \infty$ für alle k gegen Null. Hieraus folgt zusammen mit den obigen Betrachtungen und (2.1c, e)

$$\lim_{n, k \rightarrow \infty} \varrho(u_{n+r+k}, u_{n+k-1}) = 0.$$

Da $r = r(n, k)$ willkürliche ganze Zahlen ≥ 0 sind, ist das Cauchysche Konvergenzkriterium erfüllt. Daher existiert wegen der Vollständigkeit von R ein Grenzelement

$$(3.21) \quad u = \lim_{n \rightarrow \infty} u_n.$$

Für dieses Element u folgt wegen $\varrho(u_n, u_1) \leq \sigma$ (nach α) und (2.3), (2.1a, c, d, f) auf $\sigma - \varrho(u_n, u_1)$ angewandt) auch

$$(3.22) \quad \varrho(u, u_1) \leq \sigma.$$

Damit haben wir zunächst die Fehlerabschätzung (3.8). Es gilt jedoch auch die in vielen Fällen schärfere Abschätzung (3.11), aus der wegen $\tau_0 \geq \varrho(Tu_0, T_0 u_0)$ die Abschätzung (3.8) folgt; denn führen wir nur einen Schritt** $u_1 = T_0 u_0$ der Iteration aus, so hängt der Fehler $\varrho(u, u_1)$ offenbar nicht von den weiteren Operatoren T_n , $n \geq 1$, sondern nur noch von dem Operator T ab. Wir können daher für die Fehlerabschätzung $T_n = T$, $n = 1, 2, \dots$, annehmen und erhalten so als kleinste obere Schranke statt τ_0 die Größe $\varrho(Tu_0, u_1)$ und damit (3.11).

* Unter der Annahme der Existenz einer kleinsten oberen Schranke für die obigen beschränkten Folgen in N wäre der Beweis ein wenig einfacher mit dem Supremum. Wir setzen dies jedoch nicht voraus.

** Vgl. die 3. Anmerkung, S. 50.

δ) $u = \lim_{n \rightarrow \infty} u_n$ ist Lösung von (1.1). Denn es folgt aus

$$\begin{aligned} 0 \leq \varrho(u, Tu) &\leq \varrho(u, T_n u_n) + \varrho(T_n u_n, Tu) \\ &\leq \varrho(u, u_{n+1}) + \varrho(T_n u_n, T_n u) + \varrho(T_n u, Tu) \\ &\leq \varrho(u, u_{n+1}) + P \varrho(u_n, u) + \varrho(T_n u, Tu) \end{aligned}$$

und anschließendem Grenzübergang $n \rightarrow \infty$ nach (3.21) und (3.2), der Stetigkeit von P sowie (2.1 b, c, e) $\varrho(u, Tu) = 0$, also $u = Tu$.

ε) Gäbe es zwei Lösungen u und $v \neq u$ in G , so wäre nach (3.5)

$$0 \leq \varrho(u, v) = \varrho(Tu, Tv) \leq P \varrho(u, v) \leq \dots \leq P^n \varrho(u, v),$$

und hieraus folgt wegen der Konvergenz der Reihe (3.9) $\lim_{n \rightarrow \infty} P^n \varrho = 0$, also nach (2.1 a, e) $\varrho(u, v) = 0$ oder $u = v$ im Widerspruch zur Annahme, w.z.b.w.

Wir bringen noch folgenden *Zusatz zum Fixpunktsatz*:

Wird die Existenz einer Lösung u der Gleichung $u = Tu$ in G vorausgesetzt, so kann die Forderung (3.2) der gleichmäßigen Konvergenz der Folge $\{T_n\}$ gegen T für alle u von G wesentlich abgeschwächt werden. Der Satz gilt bereits mit Ausnahme der Eindeutigkeitsaussage, wenn man statt (3.2) nur voraussetzt, daß die Folge $\{T_n\}$ in dem Punkte $u = Tu$ der Lösung in G konvergiert:

$$\lim_{n \rightarrow \infty} \varrho(Tu, T_n u) = 0 \quad \text{mit} \quad u = Tu \in G.$$

Beweis. Der Teil α) des Beweises kann genau so geführt werden. Es liegen also alle u_n in G . Da auch u nach Voraussetzung in G liegt, gilt

$$\begin{aligned} \varrho(u_{n+1}, u) &= \varrho(T_n u_n, Tu) \leq \varrho(T_n u_n, T_n u) + \varrho(T_n u, Tu) \\ &\leq P \varrho(u_n, u) + \varrho(T_n u, Tu). \end{aligned}$$

Durch mehrmalige Anwendung dieser Abschätzung ergibt sich

$$(3.23) \quad \varrho(u_{n+k}, u) \leq P^k \varrho(u_n, u) + \sum_{v=1}^k P^{v-1} \varrho(T_{n+k-v} u, Tu).$$

Da alle u_n in dem durch (3.8) gegebenen Bereich liegen, folgt

$$\varrho(u_n, u) \leq \varrho(u_n, u_1) + \varrho(u_1, u) \leq \sigma + \varrho(u_1, u) = \tau, \quad \tau \in N.$$

Hieraus folgt wegen der Konvergenz der Neumannschen Reihe (3.9) für alle Punkte von N und nach (2.1 e)

$$(3.24) \quad \lim_{k \rightarrow \infty} P^k \varrho(u_n, u) = 0.$$

Nach (3.7) liegen die Größen $\varrho_n = \varrho(T_{n+k-v} u, Tu)$ unter einer gemeinsamen oberen Schranke: $\varrho_n \leq 2\tilde{u}$.

Da nach Voraussetzung und (2.1 b) $\lim_{n \rightarrow \infty} \varrho_n = 0$ ist, existiert wegen (2.1 g') nach Vorgabe einer positiven Zahl m eine Zahl $n_0(m)$, so daß

$$\varrho_n \leq \frac{2\tilde{u}}{m} \quad \text{für} \quad n \geq n_0(m) \quad \text{und} \quad k \geq v$$

gilt. Hieraus folgt wegen der Linearität und Positivität des Operators P die Abschätzung

$$\sum_{\nu=1}^k P^{\nu-1} \varrho(T_{n+k-\nu} u, T u) \leq \sum_{\nu=1}^k P^{\nu-1} \frac{2\tilde{\mu}}{m} \leq \frac{2}{m} \sum_{\nu=1}^{\infty} P^{\nu-1} \tilde{\mu}, \quad n \geq n_0(m).$$

Hieraus schließt man leicht, daß der 2. Summand auf der rechten Seite von (3.23) mit $n \rightarrow \infty$ gegen Null geht. Zusammen mit (3.24) ergibt sich sodann die Behauptung

$$\lim_{s \rightarrow \infty} \varrho(u_s, u) = 0.$$

Die Fehlerabschätzung folgt wie oben.

Ferner ist folgende Bemerkung für die Anwendungen häufig nützlich:

Die Voraussetzung (3.2) ist in vielen Fällen für die Anwendung des Satzes störend, z.B. in den Fällen, wo der Operator T nicht explizit bekannt ist, sondern nur ein Näherungsoperator \tilde{T} . Setzen wir $T_n = \tilde{T}$ für alle n , so ist (3.2) nicht erfüllt. In diesem Falle können wir natürlich auch keine Aussage darüber machen, ob die Folge u_n , $n=0, 1, 2, \dots$ gegen eine Lösung $u = Tu$ konvergiert. Das wird i. allg. nicht der Fall sein.

Trotzdem gilt die Fehlerabschätzung (3.14) auch dann noch ohne die Forderung (3.2), wenn wir statt dessen (3.5) fordern und wenn (3.8) mit $\tau_0 \geq \varrho(T_n u_0, T_0 u_0)$ und $\tau_0 \geq \varrho(T_n u_0, T u_0)$ erfüllt ist. Beides läßt sich in vielen Fällen zeigen, auch wenn T nicht genau bekannt ist.

Diese Folgerungen ergeben sich aus der einfachen Überlegung, daß wir, wenn wir nur n Schritte durchführen, die Operatoren T_r für $r \geq n$ uns durch T ersetzt denken dürfen, da wir sie ja gar nicht mehr benutzen. Dabei bleibt offenbar der Beweis des Satzes unverändert. Somit gilt der

Zusatz 2: Es gelten die Voraussetzungen A und B außer (3.2). Ferner gelte

$$\varrho(Tu, Tv) \leq P \varrho(u, v), \quad u, v \in G,$$

mit einem linearen, stetigen und positiven Operator P .

Liegen dann u_0 und alle u , für die

$$\varrho(u, u_1) \leq (E - P)^{-1} P \varrho(u_0, u_1) + (E - P)^{-1} \tau_0$$

mit einer Größe $\tau_0 \geq \varrho(T_n u_0, T_0 u_0)$ und $\tau_0 \geq \varrho(T u_0, T_0 u_0)$ gilt, in G und konvergiert die Reihe $\sum_{\nu=0}^{\infty} P^{\nu} \varrho$ für alle $\varrho \in N$, so existiert genau eine Lösung u der Gleichung $u = Tu$ in G und es gilt mit $u_{n+1} = T_n u_n$, $n=0, 1, \dots, s$, die Fehlerabschätzung

$$(3.25) \quad \varrho(u, u_{s+1}) \leq (E - P)^{-1} P \varrho(u_s, u_{s+1}) + (E - P)^{-1} \varrho(T u_s, u_{s+1})$$

für jedes $s \geq 0$.

Wählen wir hierin speziell $T_n = E$, $n=0, 1, 2, \dots$, so folgt wegen $u_n = u_0$ für alle n , da (3.4) erfüllt ist und die Bedingung (3.3) im Beweis des Fixpunktsatzes nur für Größen u, v aus der Menge der u_n verwendet wird, der

Zusatz 3: Gilt die allgemeine Voraussetzung A und

$$\varrho(Tu, Tv) \leq P \varrho(u, v), \quad u, v \in G,$$

mit einem linearen, stetigen und positiven Operator P , liegen ferner u_0 und alle u , für die

$$(3.26) \quad \varrho(u, u_0) \leq (E - P)^{-1} \varrho(Tu_0, u_0) = \sigma$$

gilt, in G und konvergiert schließlich die Reihe $\sum_{v=0}^{\infty} P^v \varrho$ für alle $\varrho \in N$, so existiert genau eine Lösung u der Gleichung $u = Tu$ in G , und für diese gilt die Abschätzung (3.26).

Dieser auch leicht auf andere Weise herleitbare Spezialfall ist ein in vielen Fällen sehr einfach anzuwendender Existenzsatz mit Fehlerabschätzung und Eindeutigkeitsaussage.

§ 4. Beispiele

1. Gewöhnliche nichtlineare Gleichungen. Nach unserem Fixpunktsatz ist es möglich, eine Begründung und Fehlerabschätzung für das bekannte Verfahren zu geben, daß man bei der Iteration

$$(4.1) \quad x_{n+1} = f(x_n)$$

die rechte Seite in eine Potenzreihe nach Potenzen von Funktionen $\varphi_j(x)$, $j=1, 2, \dots, s$, entwickelt und jeweils nur Glieder bis zu einer festen Ordnung berücksichtigt.

Sei etwa in einem Konvergenzgebiet $x \in B$, $|\varphi_j(x)| \leq c_j$, die Funktion $f(x)$ in eine Reihe der Form

$$(4.2) \quad f(x) = \sum_{\alpha_1 + \alpha_2 + \dots + \alpha_s = 0}^{\infty} a_{\alpha_1 \alpha_2 \dots \alpha_s}(x) [\varphi_1(x)]^{\alpha_1} [\varphi_2(x)]^{\alpha_2} \dots [\varphi_s(x)]^{\alpha_s}$$

entwickelbar und werden beim n -ten Schritt nur Potenzen bis zur Ordnung $\alpha_1 + \alpha_2 + \dots + \alpha_s \leq r_n$ berücksichtigt, so hat man statt (4.1) das Verfahren

$$(4.3) \quad \begin{aligned} x_{n+1} &= \sum_{\alpha_1 + \alpha_2 + \dots + \alpha_s \leq r_n} a_{\alpha_1 \dots \alpha_s}(x_n) [\varphi_1(x_n)]^{\alpha_1} \dots [\varphi_s(x_n)]^{\alpha_s} \\ &= f(x_n) - \sum_{\alpha_1 + \alpha_2 + \dots + \alpha_s \geq r_n + 1}^{\infty} a_{\alpha_1 \dots \alpha_s}(x_n) [\varphi_1(x_n)]^{\alpha_1} \dots [\varphi_s(x_n)]^{\alpha_s} \\ &= f_n - \Sigma_n = T_n x_n. \end{aligned}$$

Mit gewöhnlichem Zahlenabstand ($|x - y|$) hat man dann mit $Tx = f(x)$ in der Fehlerabschätzung (3.14) bzw. (3.25) für $\varrho(Tx_n, x_{n+1})$ die Größe $|\Sigma_n|$ einzusetzen.

1. Beispiel: Gegeben ist die Gleichung*

$$(4.4) \quad x^3 = 3R^3 \left(1 + \frac{r}{x} + \frac{4r^2}{3x^2} + \frac{5r^3}{3x^3} + \frac{r^3}{R^3} \right).$$

Es handelt sich um ein Problem bei der Bewegung eines Satelliten um einen Hauptplaneten. Dabei ist r der Radius des Satelliten, R im wesentlichen der des Hauptplaneten und x der Abstand beider Mittelpunkte. Hieraus ergibt sich, daß x, r und R reelle positive Zahlen sind. Daher gilt für die gesuchte Lösung ξ

* D. VAUGHAN [14].

von (4.4): $\xi \geq \sqrt[3]{3}R$. Dabei bedeutet ξ einen kritischen Abstand, unter dem der Satellit nicht mehr nur durch innere Gravitationskräfte zusammenhalten würde.

Die Lösung wurde in [14] * näherungsweise bestimmt, indem x/R nach Potenzen von r/R unter Vernachlässigung der 4. und höheren Potenzen von r/R entwickelt wurde.

Setzen wir zur Abkürzung $\sqrt[3]{3} = \alpha$, so haben wir

$$x = \alpha R \sqrt[3]{1 + \frac{r}{x} + \frac{4r^2}{3x^2} + \frac{5r^3}{3x^3} + \frac{r^3}{R^3}} = f(x).$$

Mit den obigen Bezeichnungen können wir schreiben: $\varphi_1(x) = r/x$ und $\varphi_2(x) = r/R$. Damit haben wir das Iterationsverfahren (4.3). Wählen wir etwa $x_0 = \alpha R$, $r_0 = 2$, $r_1 = 2$, $r_n = 3$ für $n \geq 2$, so ergibt sich nach leichter Rechnung

$$x_1 = \alpha R \left(1 + \frac{r}{3\alpha R} + \frac{r^2}{3\alpha^2 R^2} \right); \quad x_2 = \alpha R \left(1 + \frac{r}{3\alpha R} + \frac{2r^2}{9\alpha^2 R^2} \right)$$

und

$$x_3 = \alpha R \left(1 + \frac{r}{3\alpha R} + \frac{2r^2}{9\alpha^2 R^2} + \frac{32r^3}{81\alpha^3 R^3} \right).$$

Das Verfahren kommt zum Stillstand, $x_n = x_3$ für $n \geq 3$, wenn man sich auf die Entwicklung der Lösung bis zur 3. Potenz von r/R beschränkt. Daher bekommt die Fehlerabschätzung (3.11) bzw. (3.25), wenn wir als Abstand $\varrho(x, \xi)$ den Betrag der Differenz wählen, die einfache Gestalt

$$|\xi - x_3| \leq \frac{1}{1-P} |T_3 x_3 - f(x_3)|,$$

wobei hier P eine positive Zahl ist, z.B. das Maximum der Beträge von $f'(x)$ und $T_3'(x)$ in dem Bereich G .

Wir wählen für G den Bereich $x \geq \sqrt[3]{3}R$ und zur Illustration ** $r/R = 0,01$. Dann erhält man nach leichter Rechnung:

$$x_3 = T_3(x_3) = 1,445\,598\,49R \quad \text{und} \quad f(x_3) = 1,445\,598\,50R, \quad \text{sowie } P = 0,0024$$

und damit das Ergebnis:

$$|\xi - x_3| < \frac{1}{1-0,0024} \cdot 2 \cdot 10^{-8} < 2,005 \cdot 10^{-8}$$

oder

$$\xi = \underline{1,445\,598\,50 \pm 2,1 \cdot 10^{-8}}.$$

Offenbar sind alle Voraussetzungen von Zusatz 2 (§3) erfüllt.

Ein weiteres ganz anders geartetes Beispiel ist das folgende

2. Beispiel: Es ist die reelle positive Lösung ξ der Gleichung

$$(4.5) \quad z = 10 \cdot \zeta(z)$$

zu berechnen, wobei $\zeta(z)$ die Riemannsche Zetafunktion bedeutet:

$$(4.6) \quad \zeta(z) = \sum_{n=1}^{\infty} \frac{1}{n^z}.$$

* Dort noch ohne Konvergenznachweis und Fehlerabschätzung.

** Für kleinere Werte von r/R würde man noch eine größere Genauigkeit erhalten.

Hier ist der Operator $Tz = 10 \cdot \zeta(z)$ i. allg. nur näherungsweise bekannt je nach der Anzahl der Glieder der Reihe (4.6), die berücksichtigt werden. Diese Tatsache wird durch die gewöhnliche Iteration $z_{n+1} = Tz_n$ und die bekannten Fehlerabschätzungen nicht mehr erfaßt.

Die Reihe (4.6) konvergiert bekanntlich in dem Bereich Realteil $z = x \geq 1 + \varepsilon$, $\varepsilon > 0$, gleichmäßig und stellt dort eine stetige Funktion dar. Daher liegt wegen $10 - 10 \cdot \zeta(10) = -0,00995 < 0$ und $11 - 10 \cdot \zeta(11) = 0,99506 > 0$ eine reelle Lösung von (4.5) in dem Bereich $G: 10 \leq x \leq 11$.

Wir führen die Iteration

$$x_{k+1} = 10 \cdot \sum_{n=1}^{r_k} \frac{1}{n^{x_k}} = T_k x_k, \quad k = 0, 1, 2, \dots,$$

durch, wobei r_k eine nichtfallende Folge natürlicher Zahlen sei. Wählen wir etwa $r_0 = 2$, $r_1 = 3$ und $r_2 = r_3 = 4$, so erhalten wir mit $x_0 = 10$ die Iterationsfolge:

$$x_1 = 10,00976563$$

$$x_2 = 10,00986718$$

$$x_3 = 10,00987600$$

$$x_4 = 10,00987593.$$

Als Abstand $\varrho(x, y)$ wählen wir den gewöhnlichen Betrag. Dann erhält man etwa mittels der Abschätzungen

$$\left| \frac{d}{dx} T_k x \right| < 10 |\zeta'(x)| = 10 \sum_{n=1}^{\infty} \frac{\log n}{n^x} \leq 10 \left\{ \sum_{n=2}^4 \frac{\log n}{n^x} + \int_4^{\infty} \frac{\log t}{t^x} dt \right\} < 7 \cdot 10^{-3} = P \quad \text{für } x \geq 10$$

und

$$|T_4 x_4 - T x_4| = 10 \sum_{n=5}^{\infty} \frac{1}{n^{x_4}} \leq 10 \left\{ \frac{1}{5^{x_4}} + \int_5^{\infty} \frac{dt}{t^{x_4}} \right\} < 1,57 \cdot 10^{-6} \quad \text{für } x \geq 10,$$

sowie $|x_3 - x_4| < 1 \cdot 10^{-7}$ nach (3.11) die Fehlerschranke

$$|x_4 - \xi| < 1,6 \cdot 10^{-6}.$$

2. Systeme nichtlinearer Gleichungen. Eine weitere Anwendung des Fixpunktsatzes bieten die in der praktischen Mathematik häufig auftretenden Fälle, bei denen Gleichungen nur bis auf eine angebbare Genauigkeit ihrer Parameter vorliegen. In diesen Fällen ist der Operator T nicht (genau) bekannt. Wir wenden dann den Zusatz 2 an:

Es sei etwa das nichtlineare Gleichungssystem in impliziter Gestalt

$$(4.7) \quad \tilde{F}_i(x^1, x^2, \dots, x^n) = 0 \quad (i = 1, 2, \dots, n)$$

gegeben, wobei \tilde{F}_i andeuten soll, daß die Gleichungen nur näherungsweise vorliegen. Wir schreiben zur Abkürzung für (4.7)

$$(4.8) \quad \tilde{G}u = 0 \quad \text{mit dem Vektor } u = (x^1, x^2, \dots, x^n)$$

und verwenden zur Lösung z.B. das vereinfachte Newtonsche Verfahren*

$$(4.9) \quad u_{n+1} = u_n - (\tilde{G}'_{(u_0)})^{-1} \tilde{G} u_n = \tilde{T} u_n,$$

wobei $\tilde{G}'_{(u_0)}$ die an der Stelle u_0 genommene Funktionaldeterminante von (4.7) nach den x^j ist und $(\tilde{G}'_{(u_0)})^{-1}$ deren Kehrmatrix bedeutet, deren Existenz im Falle $\tilde{G}'_{(u)} \neq 0$ in einer Umgebung* \mathfrak{B} der Lösung u von (4.8) gesichert ist.

Wir verwenden als Abstand $\|u - v\|$ den Vektor, dessen Komponenten die absoluten Beträge von $u - v$ sind, und bezeichnen allgemein mit $\|A\|$ die Matrix, dessen Elemente die absoluten Beträge der Matrix A sind.

Die Ungenauigkeit der Ausgangsgleichung (4.8) sei z.B. durch

$$(4.10) \quad \|\tilde{G} u_0 - G u_0\| \leq \mu_0 \quad u_0 \in \mathfrak{B}$$

beschränkt.

Gilt dann

$$\|\tilde{T} u - \tilde{T} v\| \leq P \|u - v\| \quad \text{und} \quad \|T u - T v\| \leq P \|u - v\| \quad \text{in } \mathfrak{B}$$

mit einer positiven Matrix P und dem Operator

$$T u = u - (\tilde{G}'_{(u_0)})^{-1} G u$$

und kann man außerdem zeigen, daß $(E - P)^{-1}$ existiert und der Bereich

$$(4.11) \quad \|u - u_1\| \leq (E - P)^{-1} [P \|u_0 - u_1\| + \tau_0]$$

mit $\tau_0 = \|(\tilde{G}'_{(u_0)})^{-1}\| \mu_0$ in \mathfrak{B} liegt, so folgt die Fehlerabschätzung (3.25) mit $\varrho(u, u_{s+1}) = \|u - u_{s+1}\|$. Speziell schätzen wir ab

$$\begin{aligned} \|\tilde{T} u - \tilde{T} v\| &\leq \|(\tilde{G}'_{(u_0)})^{-1}\| [\tilde{G}'_{(u_0)}(u - v) - (\tilde{G} u - \tilde{G} v)] \\ &\leq \|(\tilde{G}'_{(u_0)})^{-1}\| \max_{u \in \mathfrak{B}} \|\tilde{G}'_{(u_0)} - \tilde{G}'_{(u)}\| \|u - v\| \end{aligned}$$

und entsprechend

$$\|T u - T v\| \leq \|(\tilde{G}'_{(u_0)})^{-1}\| \max_{u \in \mathfrak{B}} \|\tilde{G}'_{(u_0)} - G'_{(u)}\| \|u - v\|$$

und wählen die positive Matrix P so, daß $P \|u - v\|$ größer oder gleich den rechten Seiten dieser Gleichungen wird.

3. Beispiel: Als spezielles Beispiel soll die nichtlineare Behandlung der Matrizenigenwertaufgabe**,***

$$A x = \lambda B x$$

mit den Matrizen

$$A = \begin{pmatrix} 0,89689 & -2,01310 \\ 1,01900 & 1,49791 \end{pmatrix} \quad \text{und} \quad B = \begin{pmatrix} 0,03082 & -1,20201 \\ 1,13810 & 2,39792 \end{pmatrix},$$

* Wir setzen hier alle erforderlichen Stetigkeits- und Differenzierbarkeitseigenschaften der \tilde{F}_i voraus.

** Die nichtlineare Behandlung von Matrizenigenwertaufgaben führte m. W. H. UNGER [15] zuerst ein. Vgl. ferner L. COLLATZ [8], wo auch Verfahren höherer Ordnung zur Lösung benutzt werden.

*** Ohne Berücksichtigung der Ungenauigkeit der Ausgangsmatrizen kann man hier natürlich die Lösungen exakt angeben. Daher könnte der erste Teil der Rechnung fortgelassen werden. Er soll hier nur zur Illustration gebracht werden. Bei Matrizen mit mehr Zeilen würde man hier außerdem wegen der einfacheren Rechnung zweckmäßig im Banachraum mit Zahlenabstand rechnen.

deren Elemente mit einer Genauigkeit von $\pm 5 \cdot 10^{-6}$ gegeben sind, durchgerechnet werden.

Zur Normierung setzen wir die erste Komponente $x^{(1)}$ des Eigenvektors gleich 1. Dann ist die Unbekannte u der Vektor $u = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix}$. Mit $\tilde{G}u = Ax - \lambda Bx$ folgt

$$(4.12) \quad \tilde{G}'_{(u_0)} = \begin{pmatrix} -2,01310 + 1,20201\lambda & -0,03082 + 1,20201x^{(2)} \\ 1,49791 - 2,39792\lambda & -1,13810 - 2,39792x^{(2)} \end{pmatrix}.$$

Als Ausgangsvektor sei die bereits recht gute Näherung $u_0 = \begin{pmatrix} 0,768 \\ 0,728 \end{pmatrix}$ gegeben. Dann wird

$$(\tilde{G}'_{(u_0)})^{-1} = \begin{pmatrix} -0,824920876 & -0,247036883 \\ 0,068595906 & -0,315061731 \end{pmatrix} \quad \text{und} \quad \tilde{G}u_0 = \begin{pmatrix} 0,000440839 \\ 0,000171416 \end{pmatrix}.$$

Hieraus ergibt sich für den Zuwachs δ und die neue Näherung u_1 :

$$\delta = \begin{pmatrix} 0,000406003 \\ 0,000023767 \end{pmatrix}; \quad u_1 = \begin{pmatrix} 0,768406003 \\ 0,728023767 \end{pmatrix}.$$

Wir wählen nun für \mathfrak{B} den Bereich

$$\|u - u_1\| \leq 10^{-4} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

In \mathfrak{B} gilt nach (4.12)

$$\begin{aligned} \|\tilde{G}'_{(u_0)} - \tilde{G}'_{(u)}\| &\leq \|\tilde{G}'_{(u_0)} - \tilde{G}'_{(u_1)}\| + \begin{pmatrix} 1,20201 & 1,20201 \\ 2,39792 & 2,39792 \end{pmatrix} \cdot 10^{-4} \\ &\leq \begin{pmatrix} 1,4870 & 6,1822 \\ 2,9683 & 12,1343 \end{pmatrix} \cdot 10^{-4}. \end{aligned}$$

Wegen der Ungenauigkeit der Ausgangsmatrizen sind die Elemente von $\|\tilde{G}'_{(u_0)} - \tilde{G}'_{(u)}\|$ in \mathfrak{B} nach (4.12) um weniger als $1 \cdot 10^{-5}$ größer. Durch Multiplikation mit $\|(\tilde{G}'_{(u_0)})^{-1}\|$ ergibt sich so für den Abschätzungsoperator *

$$P = \begin{pmatrix} 2,0672 & 8,2047 \\ 1,0756 & 4,2855 \end{pmatrix} \cdot 10^{-4}.$$

Nach dem Zeilensummenkriterium konvergiert die Reihe $\sum_{v=0}^{\infty} P^v \varrho$ für alle $\varrho \in N$. Wir können hier jedoch leicht $(E - P)^{-1}$ genau berechnen. Man findet *

$$(E - P)^{-1} = \begin{pmatrix} 1,00021 & 0,00083 \\ 0,00011 & 1,00043 \end{pmatrix}.$$

Es ist $T_n = \tilde{T}$ für $n=0, 1, 2, \dots$. Daher brauchen wir nur noch $\|\tilde{T}u_0 - Tu_0\|$ nach oben abzuschätzen, um die Fehlerabschätzung (3.25) anwenden zu können. Man findet

$$\|\tilde{T}u_0 - Tu_0\| \leq 5 \cdot 10^{-6} (1 + \lambda_0) \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ x_0^{(2)} \end{pmatrix} \leq \begin{pmatrix} 1,528 \\ 1,528 \end{pmatrix} \cdot 10^{-5}$$

und damit *

$$\|\tilde{T}u_0 - Tu_0\| \leq \|(\tilde{G}'_{(u_0)})^{-1}\| \|\tilde{G}u_0 - Gu_0\| \leq \begin{pmatrix} 1,6380 \\ 0,5863 \end{pmatrix} \cdot 10^{-5},$$

* Dabei wurde stets nach oben abgerundet.

und schließlich ergibt (4.11) die Fehlerabschätzung

$$\|u - u_1\| \leq \begin{pmatrix} 0,104 \\ 0,055 \end{pmatrix} \cdot 10^{-6} + \begin{pmatrix} 1,6389 \\ 0,5898 \end{pmatrix} \cdot 10^{-5} = \begin{pmatrix} 1,6493 \\ 0,5923 \end{pmatrix} \cdot 10^{-5}.$$

Dabei ist der erste Summand der Summe eine Abschätzung des Fehlers, wenn die Matrizen A und B genau gegeben sind. In diesem Falle hätte man den Bereich \mathfrak{B} jedoch noch kleiner wählen können und würde einen noch kleineren Fehler $\|u - u_1\|$ erhalten. Durch die Ungenauigkeit der gegebenen Matrizen wird der Fehlerbereich wesentlich vergrößert. Offenbar liegt u in \mathfrak{B} . Daher ist die Abschätzung gültig.

3. Nichtlineare Randwertaufgaben. Bei nichtlinearen Randwertaufgaben bei gewöhnlichen Differentialgleichungen hat man in vielen Fällen mit Hilfe der Greenschen Funktion, falls sie existiert, eine Möglichkeit zu einer einfachen Fehlerabschätzung zu gelangen^{*}. Dies ist häufig auch bei der Iteration mit veränderlichen Operatoren bzw. Näherungsoperatoren der zweckmäßige Weg, um die Fehlerabschätzung (3.14) durchzuführen.

Es sei etwa die spezielle Randwertaufgabe

$$(4.13) \quad L[y] = \tilde{f}(x, y), \quad a \leq x \leq b, \quad U_i[y] = c_i \quad (i = 1, 2, \dots, m)$$

gegeben, wobei L ein linearer Differentialoperator ist und \tilde{f} wieder andeuten soll, daß die Gleichung nur bis auf eine angebbare Ungenauigkeit vorliegt. Es existiere die Greensche Funktion $G(x, \xi)$ des Randwertproblems

$$L[y] = r(x), \quad U_i[y] = c_i \quad (i = 1, 2, \dots, m).$$

Dann ist (4.13) gleichwertig mit der Gleichung

$$(4.14) \quad y = \int_a^b G(x, \xi) \tilde{f}(\xi, y(\xi)) d\xi = \tilde{T}y.$$

Die Ungenauigkeit von (4.13) sei in einem Bereich \mathfrak{B} durch

$$(4.15) \quad \varrho(\tilde{f}(x, y(x)), f(x, y(x))) \leq \mu$$

beschränkt. Dann ist

$$Ty = \int_a^b G(x, \xi) f(\xi, y(\xi)) d\xi,$$

wobei $f(x, y)$ und also auch Ty nicht (genau) bekannt sind.

Wählen wir der Einfachheit halber als Abstand den Zahlenabstand

$$\varrho(u, v) = \|u(x) - v(x)\| = \max_{x \in \mathfrak{B}} |u(x) - v(x)|$$

und genügen die Funktionen \tilde{f} und f in \mathfrak{B} einer Lipschitzbedingung

$$|\tilde{f}(x, u) - \tilde{f}(x, v)| \leq L|u - v|$$

^{*} Vgl. L. COLLATZ [6], S. 180ff. und J. SCHRÖDER [10] §§ 3 und 4.

mit der gleichen Lipschitzkonstanten L , so können wir

$$(4.16) \quad P = L \max_{x \in \mathfrak{B}} \int_a^b |G(x, \xi)| d\xi$$

setzen und die Größe $\|Tu_0 - \tilde{T}u_0\|$ durch

$$(4.17) \quad \|Tu_0 - \tilde{T}u_0\| \leq \mu \cdot \max_{x \in \mathfrak{B}} \int_a^b |G(x, \xi)| d\xi$$

abschätzen.

4. *Beispiel:* Als spezielles Beispiel betrachten wir die nichtlineare Schwingungsaufgabe

$$(4.18) \quad y'' + 4y + y^3 + \gamma y^5 = \sin x, \quad |\gamma| \leq 0, 1,$$

wobei also γ nicht genau gegeben sei. Gesucht sei eine periodische Lösung $y(x)$ mit der Periode 2π der „erzwingenden Kraft“ $\sin x$.

Hier existiert offenbar für die Randwertaufgabe

$$y'' + 4y = r(x), \quad y(0) = y(2\pi), \quad y'(0) = y'(2\pi),$$

keine Greensche Funktion*. Wir betrachten daher statt dessen das Randwertproblem (4.18), $y(0) = 0$, $y'(\pi/2) = 0$. Diese Randbedingungen sind hinreichend** für die Periodizität der Lösung mit der Periode 2π . Als Greensche Funktion für die Randwertaufgabe

$$y'' + 4y = r(x), \quad y(0) = y'(\pi/2) = 0$$

findet man

$$(4.19) \quad G(x, \xi) = \begin{cases} -\frac{1}{2} \cos 2x \sin 2\xi & \text{für } \xi \leq x \\ -\frac{1}{2} \sin 2x \cos 2\xi & \text{für } \xi \geq x. \end{cases}$$

Wir iterieren nach der Vorschrift der Näherungsgleichung

$$(4.20) \quad y''_{n+1} + 4y_{n+1} = -y_n^3 + \sin x, \quad y(0) = y'(\pi/2) = 0,$$

und beginnen mit der Lösung $y_0 = \frac{1}{3} \sin x$ der linearen Gleichung. Man erhält auf diese Weise die Näherungslösungen in der Gestalt

$$(4.21) \quad y_n = \sum_{\nu=1, 3, 5, \dots}^{k_n} a_\nu \sin \nu x,$$

wobei also nur Koeffizienten mit ungeradem Index auftreten. Dies gilt auch, wenn noch das Glied $-\gamma y^5$ hinzugenommen würde***.

* Für $r(x) = \sin 2x$ ist das Problem nicht lösbar. Vgl. auch L. COLLATZ [6] S. 25.

** H. EHREMANN [18] S. 133. Dort ist auch gezeigt, daß die Differentialgleichung (4.18) im Falle $\gamma \geq 0$ unendlich viele periodische Lösungen mit der Periode 2π besitzt. Im Falle $\gamma < 0$ folgt die Existenz mindestens einer Lösung der Randwertaufgabe (4.18), $y(0) = y'(\pi/2) = 0$ und damit mindestens einer periodischen Lösung aus einem Existenzsatz von H. EPHESER [16] S. 451f.

*** Für y_0 ist dies richtig, und setzt man (4.21) in die rechte Seite $f(y_n(x), x) = g(x)$ von (4.20) ein, so gilt wegen der Symmetrie von $f(y, x)$ $g(-x) = -g(x)$ und $g(x + \pi) = -g(x)$. Die Fourierreentwicklung von $g(x)$ enthält daher nur Glieder mit $\sin(2\nu + 1)x$, $\nu = 0, 1, 2, \dots$. Siehe R. ZURMÜHL [17] S. 313f. Diese bleiben bis auf die Koeffizienten a_ν bei der Integration (4.20) erhalten.

Man erhält mit $y_0 = \frac{1}{3} \sin x$

$$y_1 = \frac{35}{108} \sin x - \frac{1}{540} \sin 3x = 0,324074074 \sin x - 0,001851852 \sin 3x,$$

$$\begin{aligned} y_2 &= \frac{4091233}{12597120} \sin x - \frac{692891}{393660000} \sin 3x + \frac{11}{1574640} \sin 5x - \\ &\quad - \frac{7}{377913600} \sin 7x + \frac{1}{48498912000} \sin 9x \\ &= 0,324775266 \sin x - 0,001760125 \sin 3x + 0,6986 \cdot 10^{-5} \sin 5x - \\ &\quad - 1,85 \cdot 10^{-8} \sin 7x + 0,2 \cdot 10^{-10} \sin 9x. \end{aligned}$$

Zur Durchführung der Fehlerabschätzung geben wir den Bereich

$$\mathfrak{B}: \quad 0 \leq x \leq \frac{\pi}{2}, \quad \|u - y_2\| \leq 1 \cdot 10^{-3}$$

vor. Man findet $\max_{\mathfrak{B}} |u| < 0,3276 = u_m$ und für die Lipschitzkonstante L das Maximum der Ableitungen von \tilde{f} bzw. f nach u in \mathfrak{B} :

$$3u_m^2 + |\gamma| \cdot 5u_m^4 < 0,328 = L.$$

Weiter ergibt sich für die Greensche Funktion (4.19) nach leichter Rechnung die Abschätzung

$$\int_0^{\pi/2} |G(x, \xi)| d\xi \leq \frac{1}{4} \int_0^{\pi/2} [|\sin 2(x + \xi)| + |\sin 2(x - \xi)|] d\xi = \frac{1}{2}.$$

Damit wird $P = 0,164$ und mit $\|y_1 - y_2\| < 0,0008$ sowie

$$\|f(x, u) - \tilde{f}(x, u)\| = \|\gamma u^5\| \leq 0,1 \cdot u_m^5 < 0,0003774$$

ergibt (3.25) für die periodische Lösung $y(x)$ der Schwingungsgleichung (4.18) den Fehler

$$\|y - y_2\| = \max_{0 \leq x \leq \pi/2} |y(x) - y_2(x)| \leq 0,000157 + 0,000452 = 0,000609.$$

Es liegt y in \mathfrak{B} . Daher ist die Abschätzung gültig, und in \mathfrak{B} existiert genau diese eine periodische Lösung $y(x)$ bei beliebigem $|\gamma| \leq 0,1$.

Literatur

- [1] KUREPA, G.: Tableaux ramifiés d'ensembles. Espaces pseudo-distanciés. C. R. Acad. Sci. Paris **198**, 1563–1565 (1934).
- [2] KANTOROVITCH, L.: The method of successive approximation for functional equations. Acta Math. **71**, 63–97 (1939).
- [3] WEISSINGER, J.: Über das Iterationsverfahren. Z. angew. Math. Mech. **31**, 245–246 (1951).
- [4] WEISSINGER, J.: Zur Theorie und Anwendung des Iterationsverfahrens. Math. Nachr. **8**, 193–212 (1952).
- [5] COLLATZ, L.: Fehlerabschätzungen zum Iterationsverfahren bei linearen und nichtlinearen Randwertaufgaben. Z. angew. Math. Mech. **33**, 116–127 (1953).
- [6] COLLATZ, L.: Numerische Behandlung von Differentialgleichungen, 2. Aufl. Berlin 1955.

- [7] COLLATZ, L.: Einige funktionalanalytische Methoden bei der numerischen Behandlung von Differentialgleichungen. Vortrag Jahrestagung der Gesellschaft für Angew. Math. Mech., April 1958. Z. angew. Math. Mech. **38**, 264—267 (1958).
- [8] COLLATZ, L.: Näherungsverfahren höherer Ordnung für Gleichungen in Banachräumen. Arch. Rational Mech. Anal. **2**, 66—75 (1958).
- [9] SCHRÖDER, J.: Das Iterationsverfahren bei allgemeinerem Abstandsbegriff. Math. Z. **66**, 111—116 (1956).
- [10] SCHRÖDER, J.: Neue Fehlerabschätzungen für verschiedene Iterationsverfahren. Z. angew. Math. Mech. **36**, 168—181 (1956).
- [11] SCHRÖDER, J.: Nichtlineare Majoranten beim Verfahren der schrittweisen Näherung. Arch. Math. **7**, 471—484 (1956).
- [12] SCHRÖDER, J.: Über das Newtonsche Verfahren. Arch. Rational Mech. Anal. **1**, 154—180 (1957).
- [13] CARATHÉODORY, C.: Variationsrechnung und partielle Differentialgleichungen erster Ordnung, Bd. I, 2. Aufl., herausgeg. von E. HÖLDER. Leipzig 1956.
- [14] VAUGHAN, D.: On the form of satellites revolving at small distances from their primaries. Phil. Mag. **20** (1860).
- [15] UNGER, H.: Nichtlineare Behandlung von Eigenwertaufgaben. Z. angew. Math. Mech. **30**, 281—282 (1950).
- [16] EPHESER, H.: Über die Existenz der Lösungen von Randwertaufgaben mit gewöhnlichen, nichtlinearen Differentialgleichungen zweiter Ordnung. Math. Z. **61**, 435—454 (1955).
- [17] ZURMÜHL, R.: Praktische Mathematik, 2. Aufl. Berlin 1957.
- [18] EH RMANN, H.: Nachweis periodischer Lösungen bei gewissen nichtlinearen Schwingungsdifferentialgleichungen. Arch. Rational Mech. Anal. **1**, 124—138 (1957).

Mathematisches Institut A
der Technischen Hochschule
Stuttgart

(Eingegangen am 13. Juni 1959)

*Konstruktion und Durchführung von Iterationsverfahren höherer Ordnung**

HANS EH RMANN

Vorgelegt von L. COLLATZ

§ 1. Einleitung

Das Problem der Aufstellung von Iterationsverfahren für Gleichungen mit einer Zahl als Unbekannte, die mit einer vorgegebenen Ordnung konvergieren, wurde von zahlreichen Mathematikern immer wieder erneut in Angriff genommen¹. Seine Lösung geht jedoch bereits auf LEONHARD EULER zurück². Eine geschlossene vollständige Darstellung des unter gewissen Differenzierbarkeitsvoraussetzungen allgemeinsten Iterationsverfahrens von gegebener ganzzahliger Ordnung bringt E. SCHRÖDER 1870 [1]. Wenn hier trotzdem noch einmal dieser Problemkreis angeschnitten wird, obwohl das Problem der Aufstellung von Iterationsverfahren bestimmter Ordnung theoretisch vollständig gelöst ist, so geschieht dies, weil die meisten dieser Verfahren für die praktische Durchführung zu umständlich und unübersichtlich sind, so daß man sich bisher in der Praxis im allgemeinen auf die mehrfache Anwendung einfacher Verfahren mit kleinerer Konvergenzordnung beschränkte. Das Fehlen einer für die Rechnung handlichen Form hat offenbar auch dazu beigetragen, daß die Verfahren höherer Ordnung immer wieder in Vergessenheit geraten sind. Außerdem besteht meines Wissens trotz verschiedener Auszählungen der Rechenschritte³ noch keine einfache Entscheidungsmöglichkeit, welche Verfahren in speziellen Fällen am günstigsten sind.

So wird beispielsweise in der praktischen Mathematik zur Bestimmung der Wurzeln algebraischer Gleichungen unter den hier betrachteten Iterationsverfahren in den meisten Fällen das Newtonsche Verfahren, also ein Verfahren 2. Ordnung, bevorzugt. In der Tat erweist es sich für algebraische Gleichungen bis zum 5. Grade den Verfahren höherer Ordnung überlegen. Dagegen ist für algebraische Gleichungen vom 6. Grade an das Verfahren 3. Ordnung in der in § 3 beschriebenen Form in Bezug auf die Rechenarbeit im allgemeinen allen anderen Verfahren der Ordnung $k > 1$ vorzuziehen, was im Folgenden noch genauer erklärt wird.

* Diese Ausführungen sind ebenfalls wie eine frühere Veröffentlichung „Iterationsverfahren mit veränderlichen Operatoren“ meiner Habilitationsschrift entnommen. Auch hier möchte ich den Professoren L. COLLATZ, H. KÖNIG und G. SCHULZ für ihre freundliche Unterstützung herzlich danken.

¹ Vgl. u. a. das Literaturverzeichnis am Schluß dieser Arbeit.

² Siehe E. BODEWIG [2].

³ Siehe z. B. R. LUDWIG [3].

§ 2 ist ebenfalls nur als eine Einführung gedacht und bringt neben einigen Voraussetzungen und einfachen Folgerungen einige schon bekannte Ergebnisse (vgl. [1]), die hier nur etwas präzisiert werden. Die Ausführungen verfolgen im wesentlichen den Zweck, darauf hinzuweisen, daß die immer wieder unabhängig voneinander aufgestellten speziellen Iterationsverfahren höherer Ordnung sich unmittelbar aus allgemeineren Untersuchungen ergeben, die bereits E. SCHRÖDER durchgeführt hat, und daß die Aufstellung neuer solcher Verfahren keinerlei Schwierigkeiten macht.

In § 3 wird trotzdem eine neue Methode zur Aufstellung von speziellen Verfahren jeder ganzzahligen Ordnung $k > 1$ angegeben, weil es gelingt, diese speziellen Verfahren in eine Form zu bringen, die einerseits für jede Ordnung $k \geq 3$ die praktische Durchführung gegenüber den bisherigen Methoden wesentlich vereinfacht¹, andererseits eine unmittelbare Abschätzung der Rechenarbeit gestattet. Überhaupt scheint bei diesen Verfahren die Rechenarbeit (in bezug auf Multiplikationen und Divisionen) die untere Grenze zu bilden von allen entsprechenden Verfahren derselben Ordnung, wie mehrere Vergleiche zeigten, ich jedoch leider nicht allgemein beweisen konnte.

In § 4 werden einige Sätze über die Rechenarbeit und die Wahl der Ordnung für die in § 3 aufgestellten Verfahren bewiesen, die eine Entscheidung gestatten, welche Ordnung bei speziellen Gleichungen die zweckmäßige ist.

§ 5 bringt Fehlerabschätzungen für die betreffenden Iterationsverfahren und § 6 schließlich einige numerische Beispiele.

Es sei noch erwähnt, daß sich die aufgestellten Verfahren höherer Ordnung besonders gut für eine Übertragung auf allgemeinere Gleichungen z. B. im Banach-Raum eignen, da jeweils nur *eine* Division $[(F'(x_n))^{-1}]$ erforderlich ist, die dann der allgemeinen Lösung einer linearen Gleichung $[T'_n u = v$ z. B. mit der Fréchet'schen Ableitung $T'_n]$ entspricht. Dies soll an anderer Stelle gezeigt werden.

§ 2. Voraussetzungen und einfache Folgerungen

Zur Bestimmung einer Lösung² ξ der Gleichung

$$(2.1) \quad F(x) = 0$$

betrachten wir Iterationsverfahren der Gestalt³

$$(2.2) \quad x_{n+1} = f(x_n) = G(x_n, F(x_n), F'(x_n), \dots, F^{(s)}(x_n)), \quad s \geq 1, \quad n = 0, 1, 2, \dots$$

Ein solches Verfahren konvergiert für alle Anfangswerte x_0 einer Umgebung von ξ mindestens mit der Ordnung⁴ $k > 1$ gegen ξ , wenn

$$(2.3) \quad f(x) - \xi = O(|x - \xi|^k), \quad k > 1$$

¹ Bei Rechnung mit elektronischen Maschinen gestatten sie u. a. wegen der einheitlichen Gestalt eine einfache geschlossene Programmierung für jedes $k \geq 2$.

² Die Existenz von ξ setzen wir zunächst voraus.

³ Wir schließen damit alle Verfahren aus, bei denen im Ausdruck für $f(x)$ iterierte Funktionen wie z. B. $F(F(x))$ u. a. vorkommen.

⁴ Der Fall $k = 1$ macht weitere Fallunterscheidungen nötig. Wir beschränken uns hier auf überlineare Konvergenz.

gilt. Diese Bedingung ist erfüllt, wenn $f(x)$ in einer Umgebung von ξ Ableitungen bis zur k -ten Ordnung besitzt und die Gleichungen

$$(2.4) \quad f(\xi) = \xi, \quad f'(\xi) = f''(\xi) = \dots = f^{(k-1)}(\xi) = 0, \quad k > 1,$$

erfüllt sind. Ist außerdem $f^{(k)}(\xi) \neq 0$, so hat das Verfahren genau die Ordnung¹ k .

Die folgenden einfachen Sätze sind im wesentlichen bekannt:

Satz 1. *Ergibt $f(x)$ in (2.2) ein Iterationsverfahren von der Ordnung $k > 1$, so hat das Iterationsverfahren*

$$x_{n+1} = f_r(x_n), \quad n = 0, 1, 2, \dots,$$

wobei $f_r(x)$ die „ r -fach iterierte“ Funktion² ist, die Ordnung k^r .

Betrachtet man also $r > 1$ Schritte des Verfahrens (2.2) als einen einzigen, so erhöht sich damit die Ordnung $k > 1$ auf die r -te Potenz.

Satz 2. *Für die Funktionen $f(x)$ und $g(x)$ gelte*

$$f(x) - \xi = O(|x - \xi|^{k_1}) \quad \text{und} \quad g(x) - \xi = O(|x - \xi|^{k_2})$$

mit $k_1 > 1$ und $k_2 \geq 1$.

Dann konvergieren die Iterationsverfahren

$$x_{n+1} = f(g(x_n)) \quad \text{und} \quad x_{n+1} = g(f(x_n)), \quad n = 0, 1, 2, \dots,$$

für alle x_0 einer Umgebung von ξ mindestens mit der Ordnung $k_1 \cdot k_2$.

Satz 3. *Ergibt $f(x) = s_k(x)$ in (2.2) ein Verfahren ($k > 1$)-ter Ordnung gegen ξ , so erhält man mit $f^*(x) = s_k(x) + g(x)$ das allgemeinste³ Verfahren k -ter Ordnung, wenn $g(x)$ eine (in einer Umgebung von ξ eindeutig definierte) Funktion ist, die lediglich der Bedingung*

$$(2.5) \quad g(x) = O(|x - \xi|^k)$$

unterworfen ist.

Satz 4. *Ergibt $f(x) = s_{k_1}(x)$ ein Iterationsverfahren, das mit der Ordnung $k_1 > k > 1$ gegen ξ konvergiert, so ist*

$$\begin{aligned} \text{im Falle } F'(\xi) \neq 0 \quad x_{n+1} &= s_{k_1}(x_n) + |F(x_n)|^k \\ \text{im Falle } F'(\xi) = 0 \quad x_{n+1} &= s_{k_1}(x_n) + \left| \frac{F(x_n)}{F'(x_n)} \right|^k \end{aligned}$$

ein Verfahren von genau k -ter Ordnung⁴.

Beweis. Sowohl $F(x)$ im Falle $F'(\xi) \neq 0$ als auch $\frac{F(x)}{F'(x)}$ haben an der Stelle ξ eine einfache Nullstelle. Daher ist

$$|F(x)|^k \quad \text{bzw.} \quad \left| \frac{F(x)}{F'(x)} \right|^k \quad \text{gleich} \quad O(|x - \xi|^k),$$

¹ Diese Ordnungsdefinition geht auf E. SCHRÖDER zurück.

² $f_1(x) = f(x)$, $f_2(x) = f(f(x))$, ..., $f_r(x) = f(f_{r-1}(x))$: „ r -fach iterierte Funktion“.

³ Im Sinne von (2.3).

⁴ Es ist also keine wesentliche Einschränkung, wenn wir uns im Folgenden auf die Konstruktion von Iterationsverfahren mit ganzzahliger Ordnung $k > 1$ beschränken.

aber nicht $o(|x - \xi|^k)$. Zusammen mit

$$s_{k_1}(x) - \xi = O(|x - \xi|^{k_1}) = o(|x - \xi|^k)$$

folgt hieraus die Behauptung.

Haben wir bereits ein Verfahren k -ter Ordnung, $f(x) = s_k(x)$, so erhält man nach Satz 3 weitere solche Verfahren, wenn man zu $s_k(x)$ Funktionen $g(x)$ addiert, die (2.5) genügen, z. B. $g(x) = (x - \xi)^k$. Da aber ξ im allgemeinen nicht bekannt ist, müssen wir $g(x)$ durch die Funktion $F(x)$ und ihre Ableitungen ausdrücken, z. B. $g(x) = [F(x)]^k$ oder allgemeiner $g(x) = [F(x)]^k \varphi(x)$ mit einer Funktion $\varphi(x)$, die für $x \rightarrow \xi$ beschränkt bleibt. Unter gewissen einfachen Voraussetzungen über die Funktion $F(x)$ erhält man auf diese Weise auch *alle* Verfahren, die (mindestens) mit der Ordnung k gegen ξ konvergieren. Es gilt der

Satz 5. Die Funktion $F(x)$ besitze in ξ eine Nullstelle, sei in einer Umgebung von ξ differenzierbar, und es sei $F'(\xi) \neq 0$.

Ergibt dann (2.2) mit $f(x) = s_k(x)$, $k > 1$, ein Iterationsverfahren k -ter Ordnung, das gegen ξ konvergiert, so erhält man mit

$$(2.6) \quad f^*(x) = s_k(x) + [F(x)]^k \varphi(x)$$

mit einer willkürlichen Funktion $\varphi(x)$, die noch von $F(x)$ und ihren Ableitungen abhängen kann und die für $x \rightarrow \xi$ beschränkt bleibt, das allgemeinste¹ Verfahren k -ter Ordnung, bei dem im Ausdruck für $f(x)$ die unbekannte Größe ξ nicht auftritt.

Beispiel. Ersetzt man z. B. beim Newtonschen Verfahren $f(x) = x - \frac{F}{F'}$, das bekanntlich die Ordnung 2 nur bei einfacher Nullstelle ξ von $F(x)$ [$F'(\xi) \neq 0$] hat, die Funktion F durch die Funktion $F_1 = \frac{F}{F'}$, so erhält man das Verfahren $f_1(x) = x - \frac{FF'}{F'^2 - FF''}$, das auch bei mehrfacher Nullstelle von $F(x)$ die Ordnung 2 hat. Aus Satz 5 folgt, daß sich die Funktionen $f(x)$ und $f_1(x)$ nur durch einen Summanden $F^2 \cdot \varphi(x)$ unterscheiden, wobei $\varphi(x)$ im Falle einer einfachen Nullstelle ξ von F für $x \rightarrow \xi$ beschränkt bleibt, im Falle einer mehrfachen Nullstelle aber nicht beschränkt bleiben kann². Man erhält in der Tat eine solche Funktion:

$$\varphi(x) = \frac{F''}{F'(F'^2 - FF'')}.$$

Beweis von Satz 5. a) Es ist nach dem Mittelwertsatz

$$(2.7) \quad F(x) = (x - \xi) F'(\tilde{x}), \quad \tilde{x} \text{ Zwischenwert,}$$

wobei $F'(\tilde{x})$ nach Voraussetzung in einer Umgebung von ξ beschränkt und $\neq 0$ ist. Hieraus folgt

$$f^*(x) - \xi = s_k(x) - \xi + [F(x)]^k \varphi(x) = O(|x - \xi|^k) + O(|x - \xi|^k) = O(|x - \xi|^k).$$

Das Verfahren $f^*(x)$ konvergiert also mit der Ordnung $k > 1$.

b) Andererseits muß nach Satz 3 die Funktion $f^*(x)$ die Gestalt

$$f^*(x) = s_k(x) + g(x) \quad \text{mit} \quad g(x) = O(|x - \xi|^k)$$

¹ Es werden wieder nur solche Iterationsverfahren betrachtet, für die (2.3) gilt.

² Das F_1 nur eine einfache Nullstelle hat, ist Satz 5 hier anwendbar.

haben. Es muß also der Quotient

$$\frac{|g(x)|}{|x-\xi|^k} \quad \text{für } x \rightarrow \xi$$

beschränkt bleiben. Wir können daher $g(x)$ in der Form

$$(2.8) \quad g(x) = (x - \xi)^k \psi(x)$$

mit einer für $x \rightarrow \xi$ beschränkten, aber sonst willkürlichen Funktion $\psi(x)$ schreiben. Wegen $F'(\tilde{x}) \neq 0$ können wir

$$(2.9) \quad \varphi(x) = \frac{\psi(x)}{[F'(\tilde{x})]^k}$$

setzen, ohne die willkürliche Wahl von $\varphi(x)$ einerseits bzw. von $\psi(x)$ andererseits einzuschränken¹. Mit $\psi(x)$ bleibt auch $\varphi(x)$ für $x \rightarrow \xi$ beschränkt und umgekehrt. Aus (2.7), (2.8) und (2.9) folgt für $g(x)$ die Form

$$g(x) = [F(x)]^k \varphi(x), \quad \text{w.z.b.w.}$$

Zusatz. Ist $F'(\xi) = 0$, hat also $F(x)$ an der Stelle ξ eine mehrfache z. B. p -fache Nullstelle ($F'(\xi) = F''(\xi) = \dots = F^{(p-1)}(\xi) = 0$, $F^{(p)}(\xi) \neq 0$), so bleibt der Satz gültig und der Beweis unverändert, wenn man überall $F(x)$ durch $\frac{F(x)}{F'(x)}$ oder durch $F^{(p-1)}(x)$ ersetzt.

Satz 4 und Satz 5 zusammen zeigen, daß man *alle* Verfahren, die mit einer vorgegebenen Ordnung $k > 1$ gegen eine Lösung ξ von (2.1) konvergieren und für die (2.3) gilt, erfaßt, wenn man *ein* solches Verfahren für jedes ganzzahlige $k > 1$ angeben kann. Dieses Problem ist, wie schon bemerkt, bereits vollständig gelöst. In geschlossener Form geben E. SCHRÖDER [1] und später E. BODEWIG [2] für ein Verfahren k -ter Ordnung ($k \geq 2$) im Falle $F'(\xi) \neq 0$ die Formel an

$$(2.10) \quad f(x) = x + \sum_{v=1}^{k-1} (-1)^v \frac{F^v}{v!} \left(\frac{1}{F'} \frac{d}{dx} \right)^{v-1} \frac{1}{F'} \quad \text{E. SCHRÖDER 1869,}$$

wobei $\left(\frac{1}{F'} \frac{d}{dx} \right)^r$ bedeutet, daß der Operator $\frac{1}{F'} \frac{d}{dx}$ r -mal anzuwenden ist, also z. B. $\left(\frac{1}{F'} \frac{d}{dx} \right)^3 g(x) = \frac{1}{F'(x)} \frac{d}{dx} \left[\frac{1}{F'(x)} \frac{d}{dx} \left(\frac{1}{F'(x)} g'(x) \right) \right]$.

Bei einer mehrfachen Nullstelle ξ ($F'(\xi) = 0$) kann man auch hier in (2.10) überall $F(x)$ durch $H(x) = \frac{F(x)}{F'(x)}$ ersetzen, um die Ordnung k zu sichern.

Da die Formeln mit $F(x)$ statt $H(x)$ im allgemeinen jedoch wesentlich einfacher sind, setzen wir für das Folgende $F'(\xi) \neq 0$ voraus.

Für die praktische Durchführung der Verfahren werden die Ausdrücke (2.10) aber sehr bald zu unübersichtlich². Es ist daher für die Praxis von großer

¹ Man beachte, daß $F'(\tilde{x})$ bei festem ξ eine Funktion von x allein ist.

² Zum Beispiel erhält man für $k = 4$ aus (2.10) den Ausdruck

$$f(x) = x - \frac{F}{F'} - \frac{F^2 F''}{2F'^3} - \frac{F^3}{6F'^4} \left(3 \frac{F''^2}{F'} - F''' \right).$$

Eine zweimalige Durchführung des Newtonschen Verfahrens erscheint hier im allgemeinen vorteilhafter.

Bedeutung, Iterationsverfahren anzugeben, die mit höherer Ordnung ($k > 2$) konvergieren und deren Durchführung praktisch noch sinnvoll ist. Ferner wird eine einfache Entscheidungsmöglichkeit verlangt, welche Ordnung k in speziellen Fällen die zweckmäßige ist.

Mit diesen Problemen werden wir uns in den nächsten beiden Paragraphen beschäftigen.

§ 3. Aufstellung und Durchführung der Verfahren

Es gilt der

Hilfssatz 1. Die nach beiden Argumenten genügend oft differenzierbaren Funktionen $g_k(x, \xi)$ und $g_{k+1}(x, \xi)$ ergeben gegen ξ konvergente Iterationsverfahren

$$x_{n+1} = g_k(x_n, \xi) \quad \text{von } k\text{-ter Ordnung, } k > 1, \text{ ganz,}$$

$$\text{bzw. } x_{n+1} = g_{k+1}(x_n, \xi) \quad \text{von } (k+1)\text{-ter Ordnung.}$$

Ferner gelte

$$(3.1) \quad \left. \frac{\partial g_{k+1}(x, \xi)}{\partial \xi} \right|_{\xi = x} = 0.$$

Dann konvergiert das Iterationsverfahren

$$x_{n+1} = g_{k+1}(x_n, g_k(x_n, \xi)), \quad n = 0, 1, 2, \dots,$$

(mindestens) mit $(k+1)$ -ter Ordnung gegen ξ .

Der einfache Beweis ergibt sich unmittelbar durch Nachprüfen der Gln. (2.4) mit $f(x) = g_{k+1}(x, g_k(x, \xi))$.

Durch Entwicklung von $F(\xi)$ an einer Nachbarstelle x von ξ in eine Taylor-Reihe ergibt sich

$$(3.2) \quad 0 = F(\xi) = F(x) + (\xi - x) F'(x) + \frac{(\xi - x)^2}{2!} F''(x) + \dots$$

Wir setzen zunächst die Konvergenz dieser Reihe in einer Umgebung U von ξ voraus. Ferner sei in U $F'(x) \neq 0$. Dann können wir (3.2) in der Form schreiben:

$$(3.3) \quad \xi = x - \frac{F}{F'} - \frac{1}{F'} \left\{ \frac{(\xi - x)^2}{2!} F'' + \frac{(\xi - x)^3}{3!} F''' + \dots \right\}, \quad (F' \neq 0).$$

Ersetzen wir links ξ durch x_{n+1} und auf der rechten Seite x durch x_n , so erhalten wir mit

$$(3.4) \quad x_{n+1} = x_n - \frac{F_n}{F'_n} - \frac{1}{F'_n} \left\{ \frac{(\xi - x_n)^2}{2!} F''_n + \frac{(\xi - x_n)^3}{3!} F'''_n + \dots \right\}$$

ein Verfahren, das mit dem 1. Schritt die Lösung ξ von $F(x) = 0$ liefert, aber natürlich praktisch wegen des Auftretens der Größe ξ und der unendlich vielen Glieder keinen Sinn hat. Brechen wir die Reihe (3.4) nach dem 2. Gliede $\left(-\frac{F_n}{F'_n} \right)$ ab, so haben wir wieder das Newtonsche Verfahren. Brechen wir sie nach dem Gliede mit $(\xi - x_n)^k$, $k \geq 2$, ab, so erhalten wir ein Verfahren (mindestens) $(k+1)$ -ter Ordnung. In diesem Falle genügt es, wenn $F(x)$ $(k+1)$ -mal nach x differenzierbar ist:

Satz 6. Es sei $F(x)$ in einer Umgebung U einer Nullstelle ξ $(k+1)$ -mal differenzierbar, $k \geq 1$. Ferner sei dort $F'(x) \neq 0$.

Dann ist $x_{n+1} = g_{k+1}(x_n, \xi)$ mit

$$(3.5) \quad g_{k+1}(x, \xi) = x - \frac{F(x)}{F'(x)} - \frac{1}{F'(x)} \left\{ \frac{(\xi - x)^2}{2!} F''(x) + \dots + \frac{(\xi - x)^k}{k!} F^{(k)}(x) \right\}, \quad k \geq 2$$

ein Iterationsverfahren $(k+1)$ -ter Ordnung.

Beweis. Bricht man die Reihe in (3.3) nach dem $(k+1)$ -ten Gliede unter Berücksichtigung des Restgliedes ab und subtrahiert dann (3.3) von (3.5), so ergibt sich

$$g_{k+1}(x) - \xi = \frac{1}{F'(x)} \frac{(\xi - x)^{k+1}}{(k+1)!} F^{(k+1)}(\tilde{x}),$$

wobei \tilde{x} ein Zwischenwert ist, der in einem Bereich liegt, der x und ξ enthält. Daraus folgt $g_{k+1}(x) - \xi = O(|x - \xi|^{k+1})$, w.z.z.w.

Nun ist (3.5) wegen des Auftretens von ξ in dieser Form praktisch unbrauchbar. Haben wir aber bereits ein Verfahren k -ter Ordnung, $x_{n+1} = g_k(x_n)$, in dem der Wert ξ nicht vorkommt, so erhalten wir nach dem obigen Hilfssatz ein Verfahren $(k+1)$ -ter Ordnung, wenn wir in (3.5) ξ durch $g_k(x)$ ersetzen, und in diesem Verfahren tritt dann ebenfalls die Größe ξ nicht mehr auf. Auf diese Weise ist es möglich, rekursiv Verfahren beliebiger hoher Ordnung zu berechnen nach der Vorschrift¹:

$$(3.6) \quad \begin{aligned} f_2(x) &= x - \frac{F(x)}{F'(x)} \\ f_{k+1}(x) &= x - \frac{F(x)}{F'(x)} - \frac{1}{F'(x)} \left\{ \frac{(f_k(x) - x)^2}{2!} F''(x) + \dots + \frac{(f_k(x) - x)^k}{k!} F^{(k)}(x) \right\}, \\ &\quad k \geq 2. \end{aligned}$$

Man erhält z. B. auf diese Weise

$$(3.7) \quad \begin{aligned} f_3(x) &= x - \frac{F}{F'} - \frac{F^2 F''}{2 F'^3}, \\ f_4(x) &= x - \frac{F}{F'} - \frac{1}{F'} \left\{ \frac{F''}{2} \left(\frac{F}{F'} + \frac{F^2 F''}{2 F'^3} \right)^2 - \frac{F'''}{6} \left(\frac{F}{F'} + \frac{F^2 F''}{2 F'^3} \right)^3 \right\} \quad \text{usw.}^2. \end{aligned}$$

Der Vorzug von (3.6) liegt jedoch nicht in der Einfachheit und Übersichtlichkeit der Methode zur Aufstellung von Iterationsverfahren k -ter Ordnung ($k > 1$, ganz), sondern darin, daß (3.6) die Möglichkeit bietet, die Durchführung eines solchen Verfahrens k -ter Ordnung für die praktische Rechnung sehr zu erleichtern³. Der Grundgedanke hierfür besteht darin, daß es für die Anwendung eines Verfahrens k -ter Ordnung nicht erforderlich ist, das Verfahren wie etwa

¹ Vgl. auch R. ZURMÜHL [5] S. 17.

² Multipliziert man die Klammern aus, so könnten jeweils noch die Glieder mit F^k und den höheren Potenzen von F vernachlässigt werden, ohne daß sich die Ordnung k ändert, wie aus Satz 5 folgt.

³ Dies ist auch der einzige Grund, warum hier die zahlreichen Methoden zur Aufstellung von Iterationsverfahren k -ter Ordnung noch um eine weitere vermehrt werden.

in der Form (3.7) für $k=3$ oder 4 explizit anzugeben und die einzelnen Rechenoperationen an Hand dieser Formeln auszuführen, sondern man kann die obige Herleitung eines solchen Verfahrens in der Rechnung selbst vollziehen:

Satz 7. Die Funktion $F(x)$ sei in der Umgebung einer Nullstelle ξ mindestens $(k+1)$ -mal stetig differenzierbar mit $k \geq 1$, und es sei $F'(\xi) \neq 0$.

Dann ergibt folgende Vorschrift ein Iterationsverfahren, das mindestens mit der Ordnung $k+1$ gegen ξ konvergiert¹:

Es sei x_n ein Näherungswert von ξ . Dann berechne man

1. die Funktionswerte

$$\frac{1}{l!} \frac{d^l F(x_n)}{dx^l} = \frac{F_n^{(l)}}{l!} \quad \text{für } l = 0, 1, \dots, k,$$

2. die Verbesserungen

$$v_{n,2} = \frac{-F_n}{F'_n}, \quad v_{n,3} = \frac{-1}{F'_n} \left\{ F_n + v_{n,2}^2 \frac{F''_n}{2!} \right\}, \dots,$$

$$v_{n,k+1} = \frac{-1}{F'_n} \left\{ F_n + v_{n,k}^2 \frac{F''_n}{2!} + v_{n,k}^3 \frac{F'''_n}{3!} + \dots + v_{n,k}^k \frac{F_n^{(k)}}{k!} \right\}.$$

So wird der neue Näherungswert

$$(3.8) \quad x_{n+1} = x_n + v_{n,k+1}.$$

Beweis. Das Newtonsche Verfahren, $x_{n+1} = h(x_n) = x_n - \frac{F_n}{F'_n}$, ist ein Verfahren 2. Ordnung (wegen $F'(\xi) \neq 0$), denn man erhält durch Abbrechen der Reihe (3.2)

$$(3.9) \quad h(x) - \xi = \frac{-F''(\tilde{x})}{2F'(x)} (\xi - x)^2 = O(|x - \xi|^2),$$

wobei \tilde{x} ein Zwischenwert ist, der in einem Bereich liegt, der x und ξ enthält.

Ist $x_{n+1} = x_n + v_{n,r}$ ein Verfahren r -ter Ordnung für $2 \leq r \leq k$, so folgt aus Satz 6 und Hilfssatz 1, daß

$$x_{n+1} = x_n - \frac{1}{F'_n} \left\{ F_n + \frac{v_{n,r}^2}{2!} F''_n + \dots + \frac{v_{n,r}^r}{r!} F_n^{(r)} \right\}$$

ein Verfahren ist, das mindestens von $(r+1)$ -ter Ordnung konvergiert. Setzt man r der Reihe nach 2, 3, ..., k , so ergibt sich die Behauptung.

Bemerkung. Zur Durchführung des Verfahrens wird offenbar nur k -malige Differenzierbarkeit von $F(x)$ in ξ verlangt. In diesem Fall kann jedoch nicht die Ordnung $k+1$ gesichert werden².

Für die praktische Durchführung ist wesentlich, daß die Verbesserungen $v_{n,r+1}$, $r=2, 3, \dots, k$, ganze rationale Funktionen von $v_{n,r}$ sind, wobei die Koeffizienten $\frac{F^{(l)}}{l!}$ bezüglich r konstant bleiben. Nur der Grad ändert sich jeweils um Eins. Daher wird man hier zweckmäßig das Horner'sche Schema anwenden.

¹ Alle vorkommenden Verfahren konvergieren jeweils in einer Umgebung von ξ , ohne daß es hier stets ausdrücklich erwähnt wird.

² Man kann leicht ein Gegenbeispiel angeben.

Eine passende Anordnung ist die der Fig. 1. Der Kürze halber ist dort für die Verbesserungen v_r anstatt $v_{n,r}$ geschrieben. *Als wesentlicher Punkt in dem Schema ist zu beachten, daß in der Folge der Koeffizienten $c_l^{(l)} = \frac{F_n^{(l)}}{l!}$, $l=0, 1, 2, \dots, k$, die Stelle für $l=1$, also F'_n durch 0 zu ersetzen ist!*

Das Schema der Fig. 1 erweist sich besonders günstig, wenn $F(x)$ eine ganze rationale Funktion ist und man zur Berechnung der Funktionswerte $\frac{F^{(l)}}{l!}$ das gewöhnliche Horner'sche Schema benutzt; denn in diesem Fall erscheinen diese Funktionswerte gerade an den richtigen Stellen.

§ 4. Über den Rechenaufwand und die Wahl der günstigsten Ordnung

Um eine Übersicht über den Rechenaufwand der einzelnen Verfahren zu bekommen, beschränken wir uns im Folgenden auf das Auszählen der Multiplikationen und Divisionen und betrachten diese beiden Rechenoperationen in bezug auf den Rechenaufwand als gleichwertig. Ferner untersuchen wir nur Iterationsverfahren (3.8) nach Satz 7 und Fig. 1.

Aus dem Schema folgt der

Satz 8. *Gelten die Voraussetzungen des Satzes 7 und sei M die Anzahl der Multiplikationen und Divisionen zusammen, so sind zur Durchführung eines Schrittes bei einem Verfahren k -ter Ordnung mit $k \geq 2$ neben der Berechnung der k Funktionswerte*

$$F_n, F'_n, \frac{F''}{2!}, \dots, \frac{F_n^{(k-1)}}{(k-1)!}, [F_n^{(l)} = F^{(l)}(x_n)],$$

im allgemeinen Fall ($F^{(l)}(x_n) \neq 0$, $l=1, 2, \dots, k-1$)

$$(4.1) \quad M = \frac{k(k+1)}{2} - 2$$

Multiplikationen bzw. Divisionen durchzuführen¹.

Der Beweis ergibt sich unmittelbar durch Auszählen aus dem Schema (Fig. 1).

Folgerungen. Im allgemeinen Fall steigt also der Rechenaufwand bezüglich der Multiplikationen bzw. Divisionen quadratisch mit der Ordnung k . Erhöht man stattdessen die Ordnung, indem man ein Verfahren ($k^* > 1$)-ter Ordnung, z.B. das Newton'sche, mehrfach anwendet, so steigt, wie sich aus Satz 1 ergibt, der Rechenaufwand bezüglich Multiplikationen und Divisionen nur wie $\log k$. Ähnlich verhält es sich auch mit der Anzahl s der zu berechnenden Funktionswerte F_n, F'_n, F''_n, \dots , die bei Erhöhung der Ordnung k im ersten Fall, d.h. bei einem Schritt mit einem Verfahren k -ter Ordnung mit k linear ansteigt, $s=k$, dagegen bei mehrfacher Anwendung eines Verfahrens ($k^* > 1$)-ter Ordnung nur logarithmisch mit k anwächst².

¹ Dabei ist die Berechnung von $b = \frac{-1}{F'_n}$ (s. Fig. 1) nicht mitgezählt, sondern es wurde jeweils $c_l b = \frac{-c_l}{F'_n}$, $l=0, 1, \dots, k-1$, als eine Division gezählt.

² Zum Beispiel ist bei Anwendung eines Verfahrens von der Ordnung k^* nach dem Satz 1

$$s = \frac{k^* \log k}{\log k^*}.$$

					$b = -\frac{1}{F'_n}$	$x_{n+1} = x_n + v_r$
					$c_0 = F_n$	Verfahren 2. Ordnung (NEWTON)
				0		
				$c''_2 v_2$	$c'_2 v_2$	
				$c''_2 v_2 = c'_2$	$c_0 + c'_2 v_2 = c_2$	Verfahren 3. Ordnung
				$c'''_3 v_3$	$c'_3 v_3$	
				$c''_3 v_3 = c'_3$	$c_0 + c'_3 v_3 = c_3$	Verfahren 4. Ordnung
				$c'''_3 v_3 = c''_3$	$c_2 + c''_3 v_3 = c'_3$	
				$c^{IV}_k v_k$	$c'''_k v_k$	
				$c^{(k)}_k = \frac{F^{(k)}_n}{k!}$	$c''_2 + c''_k v_k = c'''_k$	Verfahren $(k+1)$ -ter Ordnung

Fig. 1. Rechenschema für ein Iterationsverfahren $(k+1)$ -ter Ordnung $(k \geq 4)$

Man erkennt hieraus unmittelbar, daß es im allgemeinen Fall ($F_n^{(l)} \neq 0$ für $l=1, 2, \dots$) gewöhnlich sinnlos ist, die Ordnung eines Verfahrens beliebig zu steigern, da die dabei erzielte Genauigkeit durch einen unverhältnismäßig hohen Rechenaufwand erkauft werden muß, der bei mehrfacher Anwendung eines Verfahrens kleinerer Ordnung (> 1) geringer ist.

Unter gewissen einfachen Voraussetzungen läßt sich eine Aussage für die günstigste Ordnung machen:

Satz 9. *Es werde die Rechenarbeit für die Berechnung der Funktionswerte $F(x)$, $F'(x)$, $\frac{F''(x)}{2!}, \dots, \frac{F^{(k-1)}(x)}{(k-1)!}$ an einer beliebigen Stelle $x = x_n$ in der Umgebung von ξ als gleich angesehen. Sie möge zur Berechnung eines dieser Funktionswerte der Rechenarbeit von a Multiplikationen bzw. Divisionen entsprechen.*

Dann ist für¹

$$a < a^* = \frac{\log \frac{16}{3}}{\log \frac{9}{8}} \approx 14,213 \quad \text{das Newtonsche Verfahren } (k=2)$$

und für

$$a > a^* \quad \text{das Verfahren 3. Ordnung}$$

unter allen Verfahren (3.8) das günstigste.

Bemerkung. Dieser Satz kann für die Fälle als Anhaltspunkt für die Wahl der Ordnung genommen werden, bei denen sich die Berechnung der Ableitungen von $F(x)$ nicht wesentlich gegenüber der von $F(x)$ selbst vereinfacht. Die Grenze a^* verschiebt sich in vielen Fällen noch etwas zugunsten des Verfahrens 3. Ordnung, also nach unten, wenn man berücksichtigt, daß die Durchführung des Verfahrens 3. Ordnung an Hand der Fig. 1 sehr einfach und übersichtlich möglich ist, was für die Berechnung von $F(x)$ und $F'(x)$ häufig nicht gilt.

Beweis von Satz 9. Die Anzahl der zu berechnenden Funktionswerte F , $\frac{F'}{1!}, \dots, \frac{F^{(r-1)}}{(r-1)!}$ bei einem Verfahren r -ter Ordnung ($r > 1$) ist r . Bei n -maliger Anwendung müssen daher nr Funktionswerte berechnet werden, was nach Voraussetzung der Rechenarbeit von $s = anr$ Multiplikationen bzw. Divisionen entspricht. Dazu kommen nach (4.1) $m = n \left[\frac{r(r+1)}{2} - 2 \right]$ solcher Rechenoperationen für die n -malige Durchführung des Verfahrens. Wir haben somit insgesamt die Arbeit

$$A = s + m = anr + \frac{n}{2} [r^2 + r - 4]$$

zu leisten für eine Gesamtordnung $k=r^n$, wie sie sich aus dem Satz 1 ergibt. Damit wird

$$(4.2) \quad A = \frac{\log k}{2 \log r} [r^2 + (2a+1)r - 4].$$

Halten wir nun die Gesamtordnung $k \geq 2$ fest², so wird die Arbeit A eine Funktion von der Ordnung r des angewandten Verfahrens allein: $A = A(r)$.

¹ Wir verstehen im Folgenden unter \log stets den natürlichen Logarithmus.

² Wir lassen dabei auch solche ganzzahligen $k \geq 2$ zu, die nicht in der Form r^n mit ganzzahligen r und n geschrieben werden können und geben in einer Ergänzung zu dem Beweis eine Begründung hierfür.

Es ist (wenn r zunächst stetig variabel angenommen wird)

$$\frac{dA}{dr} = \frac{\log k}{2} \frac{[2r^2 + (2a+1)r](\log r - 1) + r^2 + 4}{r(\log r)^2}.$$

Hieraus folgt

$$\frac{dA}{dr} > 0 \quad \text{für } r \geq 3 \quad (\text{wegen } \log r > 1 \text{ für } r \geq 3).$$

Daher ist stets

$$A(3) < A(r_1) \quad \text{für } r_1 = 4, 5, \dots$$

Das bedeutet, daß im Durchschnitt¹ die Rechenarbeit zur Erzielung einer Gesamtordnung k für *alle* Verfahren (3.8) mit der Ordnung $r > 3$ größer ist als bei Anwendung des Verfahrens der Ordnung 3. Es ist daher nur noch $r=2$ und $r=3$ zu vergleichen. Für diese Werte wird nach (4.2)

$$A(2) = \log k \frac{2a+1}{\log 2} \quad \text{und} \quad A(3) = \log k \frac{3a+4}{\log 3}.$$

Hieraus folgt

$$A(2) \leq A(3) \quad \text{für} \quad a \leq \frac{4 \log 2 - \log 3}{2 \log 3 - 3 \log 2} = \frac{\log \frac{16}{9}}{\log \frac{9}{8}} = a^*.$$

Damit ist der Satz 9 bewiesen.

Ergänzung. Unter dem „günstigsten“ Verfahren haben wir in Satz 9 das Verfahren r -ter Ordnung ($r \geq 2$) verstanden, bei dem man zur Erzielung einer vorgegebenen Gesamtordnung k mit der kleinsten Rechenarbeit, d.h. mit der geringsten Anzahl von Multiplikationen und Divisionen auskommt. Theoretisch kann dabei k eine beliebige natürliche Zahl > 1 sein. Da man bei der Rechnung aber nur stets eine ganzzahlige Anzahl n von Iterationsschritten durchführen kann und ebenfalls r eine ganze positive Zahl (≥ 2) ist, so erhält man nur ganzzahlige Gesamtordnungen der Form $k=r^n$ bei n -maliger Durchführung eines Verfahrens r -ter Ordnung. Ist nun $r_1 \neq r_2$ und kommen z.B. in der Zerlegung von r_1 und r_2 in Primfaktoren verschiedene Primzahlen vor, so ist auch stets $k_1=r_1^{n_1} \neq r_2^{n_2}=k_2$, wenn n_1 und n_2 natürliche Zahlen sind. Um zwei solche Verfahren mit den Ordnungen r_1 und r_2 zu vergleichen, haben wir im Beweis von Satz 9 auch nichtganzzahlige k zugelassen, die sich also praktisch nicht realisieren lassen. Für denjenigen, den dieser Weg über nichtrealisierbare Ordnungen k stört, sei zur Begründung hier folgender Satz angeführt:

Satz 10². Die Rechenarbeit (an Multiplikationen und Divisionen) zur Erzielung einer Gesamtordnung $k=r^n$ durch n -malige Durchführung eines Iterationsverfahrens von der Ordnung $r \geq 2$ sei gegeben durch den Ausdruck

$$(4.3) \quad A = Q(r) \log k,$$

wobei $Q(r)$ eine Funktion von r ist, die für (ganzzahlige) $r \geq 2$ definiert ist und dort positive Werte annimmt.

¹ Siehe Fußnote 2, S. 75.

² Alle auftretenden kleinen lateinischen Buchstaben bezeichnen in diesem Satz natürliche Zahlen.

Es sei ferner für ein $r_1 \geq 2$ und alle für $r_j \geq 2$, $r_j \neq r_1$,

$$Q_1 = Q(r_1) < Q(r_j) = Q_j.$$

Dann folgt aus

$$k_1 = r_1^{n_1} \leq r_j^{n_j} = k_j$$

stets

$$A_1 < A_j.$$

Dagegen existieren für jedes $j \neq 1$ (natürliche) Zahlen n_1 und n_j derart, daß für

$$k_1 = r_1^{n_1} > r_j^{n_j} = k_j$$

ebenfalls

$$A_1 < A_j$$

gilt.

Der Satz besagt, daß das Verfahren mit der Ordnung r_1 das günstigste ist in dem Sinne, daß 1. zur Erzielung derselben oder einer kleineren Gesamtordnung k eine kleinere Rechenarbeit erforderlich ist als bei allen anderen Verfahren $r_j \neq r_1$ und daß 2. auch mit geringerer Rechenarbeit eine größere Leistung, d.h. größere Gesamtordnung k möglich ist als bei Anwendung irgendeines Verfahrens $r_j \neq r_1$. Damit ist auch durch realisierbare Fälle der Begriff der günstigsten Ordnung r definiert. Da Gl. (4.2) die Gestalt (4.3) hat, ist damit unser Vorgehen im Beweis von Satz 9 hinreichend begründet.

Beweis von Satz 10. Aus $Q_1 < Q_j$ und $k_1 \leq k_j$ folgt

$$A_1 = Q_1 \log k_1 < Q_j \log k_j = A_j.$$

Um den zweiten Teil der Behauptung zu beweisen, wählen wir eine rationale Zahl $\frac{n_1}{n_j}$ mit

$$\frac{Q_j \log r_j}{Q_1 \log r_1} > \frac{n_1}{n_j} > \frac{Q_1 + Q_j}{2 Q_1} \frac{\log r_j}{\log r_1},$$

was wegen $Q_j > Q_1 > 0$ und $r_1, r_j \geq 2$ möglich ist. Hieraus folgt einerseits

$$\frac{n_1 \log r_1}{n_j \log r_j} = \frac{\log k_1}{\log k_j} > \frac{Q_1 + Q_j}{2 Q_1} > 1, \quad \text{also } k_1 > k_j, \text{ andererseits}$$

$$\frac{A_j}{A_1} = \frac{Q_j n_j \log r_j}{Q_1 n_1 \log r_1} > 1, \quad \text{w.z.b.w.}$$

Mit Satz 9 werden die Fälle nicht erfaßt, bei denen die Rechenarbeit zur Berechnung der Ableitungen $F^{(l)}(x)$ mit wachsendem l abnimmt. Dies ist insbesondere bei Polynomen der Fall. Für diese lassen sich leicht genaue Angaben über den Rechenaufwand machen. Es gilt der

Satz 11. Bei der Lösung algebraischer Gleichungen n -ten Grades durch Iteration nach (3.8) ist in bezug auf die Rechenarbeit (Anzahl der Multiplikationen und Divisionen) bis zum Grade $n=5$ das Newtonsche Verfahren und für einen Grad $n > 5$ das Verfahren 3. Ordnung unter allen Verfahren ($r \geq 2$)-ter Ordnung das günstigste.

Beweis. Eine einfache Auszählung ergibt bei einer algebraischen Gleichung n -ter Ordnung für ein Verfahren (3.8) mit der Ordnung $r \geq 2$ die Rechenarbeit^{1,2}

$$A^* = r(n+1) - 2.$$

Bei m -maliger Durchführung eines Verfahrens r -ter Ordnung haben wir somit die Rechenarbeit

$$A = mA^* = m[r(n+1) - 2]$$

und die Gesamtordnung nach Satz 4

$$k = r^m.$$

Also wird

$$A = \frac{\log k}{\log r} [r(n+1) - 2] = Q(r) \log k.$$

Nach Satz 10 ist die günstigste Ordnung r durch diejenige natürliche Zahl $r \geq 2$ gegeben, für die $Q(r)$ den kleinsten Wert hat. Es ist (r zunächst stetig variabel angenommen), wenn k festgehalten wird,

$$\frac{dA}{dr} = \log k \frac{r(n+1)(\log r - 1) + 2}{r(\log r)^2}$$

für $r \geq 3$ also wegen $\log 3 > 1$ $\frac{dA}{dr} > 0$, d. h. die relative Arbeit (zur Erreichung derselben Gesamtordnung k) nimmt zu, wenn man mit der Ordnung r über 3 hinausgeht. Als günstigste Verfahren kommen daher auch hier wiederum nur die Verfahren der Ordnungen 2 und 3 in Frage. Es ist

$$A(2) = \log k \frac{2n}{\log 2} \quad \text{und} \quad A(3) = \log k \frac{3n+1}{\log 3}.$$

Daraus folgt

$$A(2) \leq A(3), \quad \text{je nachdem} \quad n \leq \frac{\log 2}{\log \frac{3}{2}} \approx 5,885.$$

Damit ist der Satz bewiesen.

Die Sätze 9 und 11 ergeben beide, daß die Verfahren mit einer Ordnung > 3 unter den getroffenen Voraussetzungen unvorteilhaft sind. Ist $F(x)$ kein Polynom, so kann sich jedoch in einzelnen Fällen sehr wohl ein Iterationsverfahren von höherer als 3. Ordnung als das günstigste erweisen. Dies ist insbesondere dann der Fall, wenn sich die Rechenarbeit zur Berechnung der Ableitungen von $F(x)$ wesentlich gegenüber der zur Berechnung von $F(x)$ selbst verringert. Es lohnt sich in vielen Fällen, vor Beginn der Rechnung zu prüfen, welches Verfahren das günstigste ist. Hierzu bietet der folgende Satz eine einfache Möglichkeit:

Satz 12. *Entspricht die Berechnung von $F(x)$, $F'(x)$, $\frac{F''(x)}{2!}$, ..., $\frac{F^{(r-1)}(x)}{(r-1)!}$ einer Rechenarbeit von $a_0, a_1, a_2, \dots, a_{r-1}$ Multiplikationen bzw. Divisionen, so ist die Arbeit zur Erzielung einer Gesamtordnung k bei Verwendung eines Verfahrens*

¹ Siehe Anmerkung 1, S. 73.

² Im Gegensatz zu (4.1) wächst die Rechenarbeit hier nur linear mit der Ordnung k . Dabei wurde vorausgesetzt, daß die Funktionswerte $F_n^{(s)}/s!$ nach dem gewöhnlichen Horner-Schema berechnet werden.

r -ter Ordnung ($r \geq 2$) nach (3.8) und Fig. 1 gegeben durch

$$A = \frac{\log k}{2} \frac{r(r+1) + 2(a_0 + a_1 + \dots + a_{r-1}) - 4}{\log r}.$$

Als günstigstes Verfahren erhält man daher dasjenige der Ordnung $r \geq 2$ (ganz), bei dem der Quotient

$$(4.4) \quad Q = \frac{r(r+1) + 2(a_0 + a_1 + \dots + a_{r-1}) - 4}{\log r}$$

seinen kleinsten Wert annimmt.

Beweis. Bei einem Schritt mit einem Verfahren der Ordnung r sind nach Satz 8 Gl. (4.1) $M = \frac{r(r+1)-4}{2}$ Multiplikationen bzw. Divisionen durchzuführen. Dazu kommen nach Voraussetzung zur Berechnung von $\frac{F^{(l)}(x)}{l!}$ ($l=0, 1, 2, \dots, r-1$) $a_0 + a_1 + \dots + a_{r-1}$ solche Rechenoperationen. Bei n Schritten haben wir somit die Gesamtarbeit

$$A = n \frac{r(r+1) + 2(a_0 + a_1 + \dots + a_{r-1}) - 4}{2}$$

und die Gesamtordnung $k=r^n$, also $A = \frac{\log k}{2} Q$.

Alles andere folgt aus Satz 10.

Als Beispiel wollen wir den Fall betrachten, daß die Rechenarbeit zur Berechnung von $F(x)$ an einer Stelle x_n in der Umgebung von ξ a -Multiplikationen entspricht, die zur Berechnung der Ableitungen aber nur $\frac{a}{4}$ Multiplikationen. Dann ergibt sich nach leichter Rechnung als das günstigste Verfahren (3.8)¹

das Newtonsche Verfahren für $a < 5,019$

das Verfahren 3. Ordnung für $5,019 < a < 20,68$

das Verfahren 4. Ordnung für $20,68 < a < 117,16$

das Verfahren 5. Ordnung für $a > 117,16$.

§ 5. Fehlerabschätzungen²

Aus der Definition der Ordnung ergibt sich, daß für ein Verfahren k -ter Ordnung eine Konstante K_k existiert, so daß

$$(5.1) \quad \frac{|f(x) - \xi|}{|x - \xi|^k} \leq K_k$$

in der Umgebung der Lösung ξ von $x=f(x)$ bzw. $F(x)=0$ gilt. Daher muß in jeder Fehlerabschätzung für ein Iterationsverfahren k -ter Ordnung in irgendeiner Weise eine Abschätzung der Konstanten K_k enthalten sein, was im Falle von k -mal differenzierbarem $f(x)$ einer Abschätzung der k -ten Ableitung von $f(x)$

¹ Es sei noch einmal darauf hingewiesen, daß diese Abschätzungen nur für die speziellen nach Satz 7 durchgeführten Verfahren gelten.

² L. COLLATZ bringt in [4] Abschätzungen, die es ermöglichen, für sehr allgemeine Iterationsverfahren beliebiger Ordnung k Fehlerschranken aufzustellen. Hiervon ist bei den obigen spezielleren Fehlerabschätzungen, die sich nur auf die Verfahren des Satzes 7 beziehen, kein Gebrauch gemacht worden.

entspricht, die wiederum im allgemeinen von den k ersten Ableitungen von $F(x)$ abhängt. Hieraus folgt, daß so einfache Fehlerabschätzungen wie sie etwa für lineare Iterationsverfahren ($k=1$) und auch noch für quadratisch konvergente Verfahren möglich sind, für höhere Ordnungen nicht erwartet werden können¹.

Wir betrachten auch hier wieder nur die speziellen nach Satz 7 und Fig. 1 durchgeführten Iterationsverfahren k -ter Ordnung. Ferner sei wieder $F'(\xi) \neq 0$.

Nach (3.2) und (3.3) ist für $k \geq 3$

$$(5.2) \quad \xi = x_n - \frac{F_n}{F'_n} - \frac{1}{F'_n} \left\{ \sum_{r=2}^{k-1} \frac{(\xi - x_n)^r}{r!} F_n^{(r)} + R_{k-1,n} \right\}$$

mit

$$(5.3) \quad R_{k-1,n} = \int_0^1 \frac{t^{k-1}}{(k-1)!} (\xi - x_n)^k F^{(k)}(\xi + t(x_n - \xi)) dt, \quad \text{also} \\ |R_{k-1,n}| \leq \frac{|\xi - x_n|^k}{k!} |F^{(k)}(\xi + \vartheta(x_n - \xi))|, \quad 0 < \vartheta < 1.$$

Andererseits gilt für unsere Verfahren k -ter Ordnung nach Satz 7

$$(5.4) \quad x_{n+1} = x_n - \frac{F_n}{F'_n} - \frac{1}{F'_n} \left\{ \sum_{r=2}^{k-1} \frac{(\tilde{x}_{n+1} - x_n)^r}{r!} F_n^{(r)} \right\},$$

wobei

$$\tilde{x}_{n+1} = f(x_n)$$

sich aus einem solchen Verfahren $(k-1)$ -ter Ordnung berechnet, d.h. es ist

$$\tilde{x}_{n+1} - \xi = g_{k-1}(x_n) (x_n - \xi)^{k-1}$$

mit einer in einer Umgebung von ξ beschränkten Funktion $g_{k-1}(x)$. Daher gilt mit $g_{k-1} = g_{k-1}(x_n)$ und $k \geq 3$:

$$\begin{aligned} (\tilde{x}_{n+1} - x_n)^r &= [(\tilde{x}_{n+1} - \xi) + (\xi - x_n)]^r \\ &= (\xi - x_n)^r [1 - g_{k-1} \cdot (x_n - \xi)^{k-2}]^r \\ &= (\xi - x_n)^r + (\xi - x_n)^r \cdot \sum_{\nu=1}^r (-1)^\nu \binom{r}{\nu} g_{k-1}^\nu \cdot (x_n - \xi)^{\nu \cdot (k-2)}. \end{aligned}$$

Und hieraus folgt mit (5.2) und (5.4) für $k \geq 3$:

$$(5.5) \quad x_{n+1} - \xi = \frac{-1}{F'_n} \sum_{r=2}^{k-1} \left\{ \frac{F_n^{(r)}}{r!} \sum_{\nu=1}^r (-1)^{r-\nu} \binom{r}{\nu} g_{k-1}^\nu \cdot (x_n - \xi)^{\nu \cdot (k-2) + r} \right\} + \frac{R_{k-1,n}}{F'_n}.$$

Wegen $\nu \geq 1$, $k \geq 3$, $r \geq 2$ ist $\nu \cdot (k-2) + r \geq k$. Damit ergibt (5.5) und (5.3) wiederum

$$x_{n+1} - \xi = O(|x_n - \xi|^k).$$

Die Gl. (5.5) kann nun als Ausgangsgleichung für eine Fehlerabschätzung dienen:

¹ Die Schwierigkeit besteht hier in der Bestimmung einer möglichst kleinen Konstanten K_k in (5.1). Ist diese bekannt, so macht die Fehlerabschätzung keine Mühe.

Nehmen wir zunächst den wichtigsten Fall $k=3$ vorweg, so ergibt (5.5) und (5.3)

$$(5.6) \quad \begin{aligned} x_{n+1} - \xi &= \frac{-1}{F'_n} \frac{F''_n}{2} [-2g_2(x_n - \xi)^3 + g_2^2(x_n - \xi)^4] + \frac{R_{2,n}}{F'_n} \\ &= \frac{(\xi - x_n)^3}{2F'_n} \left[-2F''_n g_2 + F''_n g_2^2(x_n - \xi) + \frac{1}{3} \tilde{F}''' \right], \end{aligned}$$

wobei \tilde{F}''' bedeutet, daß die Funktion $F'''(x)$ an einer Zwischenstelle zu nehmen ist.

Aus (3.9) folgt

$$g_2 = \frac{-\tilde{F}'''}{2F'_n},$$

so daß sich

$$(5.7) \quad x_{n+1} - \xi = \frac{(\xi - x_n)^3}{2F'_n} \left[\frac{F''_n \tilde{F}'''}{F'_n} + \frac{F''_n \tilde{F}'''^2}{4F_n'^2} (x_n - \xi) + \frac{1}{3} \tilde{F}'''' \right]$$

ergibt. Aus dieser letzten Gleichung folgt nun leicht folgende Fehlerabschätzung für das Verfahren 3. Ordnung:

Satz 13. *Es existiere ein Bereich \mathfrak{B} , in dem folgendes gilt:*

$$\alpha) \quad |F'(x)| \geq m > 0 \quad \text{und} \quad |F^{(l)}(x)| \leq M_l$$

mit Konstanten m und M_l , $l=2, 3$.

$$\beta) \quad \frac{|F(x_0)|}{m} \leq c$$

mit einer Konstanten c .

$\gamma)$ Es sei

$$K c^2 < 1 \quad \text{mit} \quad K = \frac{1}{2m^3} \left[m M_2^2 + \frac{1}{4} M_2^3 c + \frac{1}{3} m^2 M_3 \right].$$

Liegen dann alle x mit

$$(5.8) \quad |x - x_0| \leq 2c$$

in \mathfrak{B} , so konvergiert das Iterationsverfahren (3.8) (Satz 7) für $k=3$ gegen die einzige Lösung ξ von $F(x)=0$ in \mathfrak{B} und es gilt die Fehlerabschätzung

$$(5.9) \quad |x_n - \xi| \leq K^{\frac{1}{2}(3^n - 1)} c^{3^n}, \quad n = 0, 1, 2, \dots$$

Beweis. Wegen $\alpha)$ und $\beta)$ existiert genau eine Nullstelle ξ von $F(x)$ in \mathfrak{B} , und es ist sogar

$$(5.10) \quad |\xi - x_0| \leq c.$$

Es gelte $|x_n - \xi| \leq c$ für $n=0, 1, \dots, r$. Dann liegen wegen

$$|x_n - x_0| \leq |x_0 - \xi| + |x_n - \xi| \leq 2c$$

diese x_n in \mathfrak{B} , und es folgt aus (5.7), $\alpha)$ und $\gamma)$

$$\begin{aligned} |x_{r+1} - \xi| &\leq |x_r - \xi|^3 \frac{1}{2m} \left[\frac{M_2^2}{m} + \frac{M_2^3}{4m^2} c + \frac{M_3}{3} \right] = K |x_r - \xi|^3 \\ &\leq K c^3 < c. \end{aligned}$$

Daher liegen alle x_n in dem Bereich $|x - \xi| \leq c$ von \mathfrak{B} und es gilt

$$(5.11) \quad |x_{n+1} - \xi| \leq K |x_n - \xi|^3, \quad n = 0, 1, 2, \dots$$

Wegen $K|x_n - \xi|^2 \leq Kc^2 < 1$ konvergiert das Iterationsverfahren, und durch mehrmalige Anwendung von (5.11) ergibt sich mit (5.10) die Fehlerabschätzung (5.9), wie man auch leicht durch vollständige Induktion nachprüft, w.z.b.w.

Wegen $F'(\xi) \neq 0$ läßt sich offenbar stets ein solcher Bereich \mathfrak{B} angeben, der die geforderten Bedingungen erfüllt, wenn nur x_0 nahe genug bei ξ liegt¹. Daher ist die Fehlerabschätzung immer durchführbar.

Für die Verfahren (3.8) der Ordnung $k > 3$ wird die Bestimmung der Konstanten K_k in (5.1) etwas mühsamer.

Aus (5.5) und (5.3) folgt wegen $\nu \cdot (k-2) + r \geq \nu + k - 1$ für $\nu \geq 1$, $r \geq 2$ und $k \geq 3$, wenn wir bereits wissen, daß $|x_n - \xi| \leq 1$ ist²,

$$(5.12) \quad \begin{aligned} |x_{n+1} - \xi| &\leq \frac{|x_n - \xi|^{k-1}}{|F'_n|} \left\{ \sum_{r=2}^{k-1} \frac{|F_n^{(r)}|}{r!} \sum_{\nu=1}^r \binom{r}{\nu} |g_{k-1}^\nu| |x_n - \xi|^\nu \right\} + \frac{|\tilde{F}^{(k)}|}{k! |F'_n|} |x_n - \xi|^k \\ &\leq \frac{|x_n - \xi|^{k-1}}{|F'_n|} [(1 + |g_{k-1}| |x_n - \xi|)^{k-1} - 1] \cdot \sum_{r=2}^{k-1} \frac{|F_n^{(r)}|}{r!} + \frac{|\tilde{F}^{(k)}|}{k! |F'_n|} |x_n - \xi|^k. \end{aligned}$$

Dabei bedeutet $\tilde{F}^{(k)}$, daß der Funktionswert von $F^{(k)}(x)$ an einer Zwischenstelle zwischen x_n und ξ zu nehmen ist. Die Summe

$$S_{k-1} = \sum_{r=2}^{k-1} \frac{|F_n^{(r)}|}{r!} = \frac{|F'_n|}{2!} + \frac{|F''_n|}{3!} + \dots + \frac{|F_n^{(k-1)}|}{(k-1)!}$$

kann unmittelbar aus dem Schema (Fig. 1) entnommen werden. Die Gl. (5.12) liefert eine Rekursionsformel zur Abschätzung der Größen $g_l = g_l(x_n)$:

$$(5.13) \quad |g_l| \leq \frac{1}{|F'_n|} \left\{ \frac{S_{l-1}}{|x_n - \xi|} [(1 + |g_{l-1}| |x_n - \xi|)^{l-1} - 1] + \frac{|\tilde{F}^{(l)}|}{l!} \right\}, \quad l \geq 3.$$

Hieraus ergeben sich mit

$$(5.14) \quad |F'(x)| \geq m > 0, \quad |F^{(l)}(x)| \leq M_l, \quad l = 2, 3, \dots, k, \quad x \in \mathfrak{B}$$

die Konstanten K_k in (5.1) rekursiv aus der Formel

$$(5.15) \quad K_l = \frac{1}{m} \left\{ \frac{S_{l-1}^*}{|x_n - \xi|} [(1 + K_{l-1} |x_n - \xi|)^{l-1} - 1] + \frac{M_l}{l!} \right\}, \quad l \geq 3,$$

wobei $S_{l-1}^* = \sum_{r=2}^{l-1} \frac{M_r}{r!}$ und $K_2 = \frac{M_2}{2m}$ ist.

Analog zu Satz 13 ergibt sich damit der

Satz 14. *Gelten in einem Bereich \mathfrak{B} die Abschätzungen (5.14), ist $\frac{|F(x_0)|}{m} \leq c < 1$, liegen ferner alle x mit $|x - x_0| \leq 2c$ in \mathfrak{B} und gilt $K_k c^{k-1} < 1$, wobei sich K_k*

¹ In diesem Falle muß jedoch die Existenz der Lösung ξ von $F(x) = 0$ vorausgesetzt werden.

² Dies ist offenbar der Fall, wenn $|F'(x)| \geq m > 0$ und $|F(x_n)| \leq cm < m$ in $|x - x_n| \leq c$ gilt.

rekursiv aus den Formeln

$$(5.16) \quad K_2 = \frac{M_2}{2m}, \quad K_l = \frac{1}{m} \left\{ \frac{S_{l-1}^*}{c} [(1 + K_{l-1}c)^{l-1} - 1] + \frac{M_l}{l!} \right\}, \quad 3 \leq l \leq k,$$

berechnet, so konvergiert das Verfahren (3.8) k -ter Ordnung gegen die einzige Lösung ξ von $F(x)=0$ in \mathfrak{B} , und es gilt die Fehlerabschätzung

$$(5.17) \quad |x_n - \xi| \leq K_k^{\frac{k^n-1}{k-1}} c^{k^n}, \quad n = 0, 1, 2, \dots$$

Der Beweis ist völlig analog dem von Satz 13. Auch hier existiert wegen $F'(\xi) \neq 0$ stets ein Bereich \mathfrak{B} , der die geforderten Bedingungen erfüllt. Wegen $c < 1$ kann man die 2. Gleichung in (5.16) auch durch die größere Rekursionsformel

$$K_l = \frac{1}{m} \left\{ S_{l-1}^* [(1 + K_{l-1})^{l-1} - 1] + \frac{M_l}{l!} \right\}, \quad 3 \leq l \leq k,$$

ersetzen.

§ 6. Beispiele

1. Beispiel. $e^z = 1 + z$, $z = x + yi \neq 0$.

α) Als Näherungswert für eine Lösung ξ ergibt sich aus der Zeichnung (Fig. 2) als Schnitt der Kurven gleichen Betrages und gleichen Winkels¹ der Wert $z_0 = 2,1 + 7,5i$.

β) Entspricht die Berechnung von $F_n = e^{z_n} - z_n - 1$ und $-1/F'_n$ der Rechenarbeit von $a = a_0 + a_1$ Multiplikationen (komplexer Zahlen) und vernachlässigt man dabei die Arbeit zur Berechnung der Zahlen² $F_n^{(s)}/s!$ ($s \geq 2$), so ergibt sich nach Satz 12 als die günstigste Ordnung r diejenige, bei der der Quotient (4.4)

$$(6.1) \quad Q = \frac{r(r+1) + 2a - 4}{\log r}$$

seinen kleinsten Wert³ annimmt. Es ist also das Verfahren der Ordnung $k=r$ ($r \geq 2$) günstiger (ungünstiger) als das Verfahren der Ordnung $k=r+1$, wenn $Q(r)$ kleiner (größer) als $Q(r+1)$ ist, woraus nach (6.1)

$$a \stackrel{<}{>} \frac{1}{2} \frac{[(r+1)(r+2) - 4] \log r - [r(r+1) - 4] \log(r+1)}{\log(r+1) - \log r} = G(r)$$

folgt. Es ergibt⁴ sich für

$r =$	2	3	4	5	6	7	8	9
$G(r) =$	4,129	11,275	23,063	39,966	62,360	90,586	124,902	165,526

¹ Beide Kurven lassen sich leicht punktweise mit Hilfe der Kreis- bzw. Geraden-schar durch den Punkt $(-1, 0)$ konstruieren. Sie sind symmetrisch zur x -Achse.

² Diese lassen sich etwa für $s \leq 10$ unmittelbar fortlaufend hinschreiben. Wesentlich zur Berechnung der Funktionswerte $F_n^{(s)}/s!$ ist hier in bezug auf die Rechenarbeit nur die Berechnung von

$$e^{z_n} = e^{x_n} (\cos y_n + i \sin y_n).$$

³ Die Basis des Logarithmus ist dabei beliebig > 1 .

⁴ In der Praxis würde natürlich hier eine viel größere Schätzung genügen.

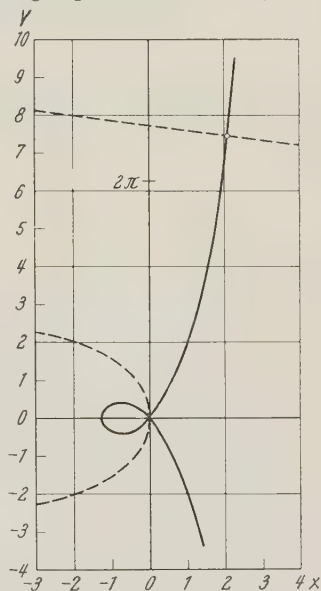


Fig. 2. Erläuterung zum 1. Beispiel

Ist also z.B. $a < 4,129$, so ist das Verfahren 2. Ordnung günstiger als das 3. Ordnung; ist $a < 11,275$, so ist das Verfahren 3. Ordnung günstiger als das 4. Ordnung usw.

Die Wahl der Ordnung hängt hier wesentlich davon ab, welche Genauigkeit erzielt werden soll. Für größere Genauigkeiten werden die üblichen Tafeln im allgemeinen nicht mehr ausreichen, und man hat die Funktionswerte (insbesondere e^{z_n}) selbst zu berechnen. Daher wird man in diesem Fall zweckmäßig ein Verfahren höherer Ordnung wählen¹.

Tabelle 1. Iterationsverfahren k -ter Ordnung $z_1 = z_0 + v_k$ zur Bestimmung

			$\frac{F''(z_0)}{2!} =$
			1,415 341 38 + 3,829 933 68 i
		$\frac{F'''(z_0)}{3!} =$	0,043 902 37 - 0,032 443 96 i
		0,471 780 46 + 1,276 644 56 i	1,459 243 75 + 3,797 489 72 i
	$\frac{F^{(4)}(z_0)}{4!} =$	0,010 975 56 - 0,008 103 16 i	0,043 467 72 - 0,032 744 91 i
	0,117 945 12 + 0,319 161 14 i	0,482 756 02 + 1,268 541 40 i	1,458 809 10 + 3,797 188 77 i
$\frac{F^{(5)}(z_0)}{5!} =$	0,002 195 05 - 0,001 620 60 i	0,010 888 33 - 0,008 169 48 i	0,043 464 82 - 0,032 740 27 i
0,023 589 02 + 0,063 832 23 i	0,120 140 17 + 0,317 540 54 i	0,482 668 79 + 1,268 475 08 i	1,458 806 20 + 3,797 193 41 i

γ) Wir führen hier die Verfahren bis zur 6. Ordnung durch (Tabelle 1). Dabei wurde 8-stellig gerechnet, um Abrundungsfehler zu vermeiden. Diese können sich hier nicht aufschaukeln, wie man sich leicht an Hand der Rechnung überzeugt. Geht man mit der Ordnung k über 6 hinaus, so ändern sich die v_k nicht mehr innerhalb der Genauigkeit von 8 Stellen nach dem Komma:

$$v_k = v_6 \pm 0,5 \cdot 10^{-8}.$$

Mit $z_1 = z_0 + v_k$ und $z_0 = 2,1 + 7,5 i$ erhalten wir somit die folgenden Ergebnisse bei einem Schritt mit dem

Verfahren 2. Ordnung (NEWTON): $z_1 = 2,088\,205\,9 + 7,462\,021\,7 i$

Verfahren 3. Ordnung: $z_1 = 2,088\,821\,5 + 7,461\,480\,1 i$

Verfahren 4. Ordnung: $z_1 = 2,088\,843\,0 + 7,461\,488\,2 i$

Verfahren 5. Ordnung: $z_1 = 2,088\,843\,0 + 7,461\,489\,2 i$

Verfahren 6. Ordnung: $z_1 = 2,088\,843\,0 + 7,461\,489\,3 i.$

¹ Wegen der einfachen geschlossenen Programmierung der Verfahren gilt dies auch bei Benutzung einer elektronischen Rechananlage.

δ) Eine grobe Schätzung läßt erwarten, daß für den Konvergenznachweis und die Fehlerabschätzung nach Satz 14 der Bereich

$$(6.2) \quad \mathfrak{B}: |z - z_0| \leq 0,1$$

die geforderten Bedingungen erfüllt. Aus (6.2) folgt mit $z = x + iy$

$$\left. \begin{array}{l} 2,0 \leq x \leq 2,2 \\ 7,4 \leq y \leq 7,6 \end{array} \right\} \text{ und daraus } |F'(z)| = |e^z - 1| \geq |e^x| - 1 = e^x - 1, \\ \text{also } |F'(z)| \geq e^2 - 1 = 6,389 = m, \quad z \in \mathfrak{B}.$$

einer Lösung der Gleichung $F(z) = e^z - z - 1 = 0$, ($z_0 = 2,1 + 7,5i$)

		$\frac{-1}{F'(z_0)} =$
$F(z_0) =$		$-0,029\,515\,26 + 0,123\,496\,52i$
0	$-0,269\,317\,25 + 0,159\,867\,37i$	$v_0 = -0,011\,794\,10 - 0,037\,978\,27i$
	$-0,005\,275\,55 - 0,003\,723\,44i$	
$0,128\,761\,58 - 0,098\,922\,84i$	$-0,274\,592\,80 + 0,156\,143\,93i$	$v_1 = -0,011\,178\,55 - 0,038\,519\,88i$
	$-0,005\,253\,22 - 0,003\,903\,42i$	
$0,129\,966\,62 - 0,098\,660\,32i$	$-0,274\,570\,47 + 0,155\,963\,95i$	$v_2 = -0,011\,156\,99 - 0,038\,511\,81i$
	$-0,005\,245\,18 - 0,003\,905\,54i$	
$0,129\,960\,69 - 0,098\,546\,58i$	$-0,274\,562\,43 + 0,155\,961\,83i$	$v_3 = -0,011\,156\,96 - 0,038\,510\,76i$
	$-0,005\,244\,96 - 0,003\,905\,28i$	
$0,129\,956\,96 - 0,098\,544\,87i$	$-0,274\,562\,21 + 0,155\,962\,09i$	$v_4 = -0,011\,157\,00 - 0,038\,510\,74i$

Weiter wird

$$\frac{|F(z_0)|}{m} = 0,0490 < 0,05 = c < 1,$$

und damit gilt $|z - z_0| \leq 2c$ in \mathfrak{B} .

Ferner ergeben sich unmittelbar die Abschätzungen

$$|F^{(l)}(z)| = |e^z| = e^x \leq e^{2,2} = 9,025\,013\,5 \leq 9,0251 = M_l, \quad l = 2, 3, \dots$$

in \mathfrak{B} , und hieraus rekursiv aus der Formel (5.16) folgende Konstanten K_k bei den Verfahren der Ordnung:

$$(6.3) \quad \begin{array}{c|c|c|c|c} k = & 2 & 3 & 4 & 5 & 6 \\ \hline K_k = & 0,706\,31 & 1,2508 & 3,8183 & 20,2544 & 648,556 \end{array}.$$

Es folgt $K_k c^{k-1} < 1$ für $k = 2, 3, \dots, 6$. Daher konvergieren diese Verfahren.

Unter Verwendung der Konstanten K_k ergeben sich jeweils nach dem n -ten Schritt aus der Fehlerformel (5.17) die Schranken für das Verfahren

- 2. Ordnung: $|z_n - \zeta| \leq 1,42 (3,532 \cdot 10^{-2})^{2^n}$
- 3. Ordnung: $|z_n - \zeta| \leq 0,895 (5,592 \cdot 10^{-2})^{3^n}$
- 4. Ordnung: $|z_n - \zeta| \leq 0,639 (7,830 \cdot 10^{-2})^{4^n}$
- 5. Ordnung: $|z_n - \zeta| \leq 0,472 (1,0608 \cdot 10^{-1})^{5^n}$
- 6. Ordnung: $|z_n - \zeta| \leq 0,274 (1,8255 \cdot 10^{-1})^{6^n}$.

Noch genauere Fehlerschranken erhält man, wenn man die Fehlerabschätzungen von Satz 14 erst *nach* dem ersten Schritt durchführt. Für den Wert $z_1 = 2,0888430 + 7,4614893i$ des Verfahrens 6. Ordnung ergibt sich dabei unter Verwendung des neuen Bereiches \mathfrak{B}_1 , $|z - z_1| \leq 10^{-7}$, ein Fehler von $|z_1 - \xi| < 0,15 \cdot 10^{-7}$. z_1 ist also noch in den letzten Stellen genau.

2. Beispiel. $F(x) \equiv x^7 + 5x^6 + 3x^5 + 2x^4 + 4x^3 + 2x^2 + 6x + 1 = 0$.

Es ist $F(-1) = -1$ und $F(-0,5) = 1,1016$. Daher liegt zwischen $x = -1$ und $x = -0,5$ eine reelle Nullstelle von $F(x)$. Wir wählen $x_0 = -0,75$.

Nach Satz 11 ist hier das Verfahren 3. Ordnung das günstigste.

Zur Veranschaulichung der Verfahren wollen wir hier trotzdem das Schema der Fig. 1 über $k=3$ hinaus weiterrechnen. Wir erhalten so die Tabelle 2 und damit das Ergebnis nach einem Schritt bei dem

Verfahren 2. Ordnung (NEWTON):	$x_1^{(2)} = -0,6765992$
Verfahren 3. Ordnung:	$x_1^{(3)} = -0,6825957$
Verfahren 4. Ordnung:	$x_1^{(4)} = -0,6807081$
Verfahren 5. Ordnung:	$x_1^{(5)} = -0,6809935$
Verfahren 6. Ordnung:	$x_1^{(6)} = -0,6809586$
Verfahren 7. Ordnung:	$x_1^{(7)} = -0,6809626$
Verfahren 8. Ordnung:	$x_1^{(8)} = -0,6809622$

Bei Weiterrechnung des Schemas ergibt sich durch die Verfahren 9. und höherer Ordnung innerhalb der Rechengenauigkeit von 7 Stellen keine Änderung mehr.

Führen wir statt dessen zweimal das (günstigste) Verfahren 3. Ordnung durch, so erhalten wir (mit $x_0 = -0,75$; $x_1 = -0,6825957$) den Wert $x_2 = -0,6809622$, also denselben wie bei einem Schritt mit dem Verfahren 8. Ordnung.

Es ist $F(-0,6809622) = 4,4 \cdot 10^{-8}$, also $|F(x_1^{(8)})| < 5 \cdot 10^{-8}$. In dem Bereich \mathfrak{B} : $-0,681 \leq x \leq -0,680$ ist $|F'(x)| > 5,8 = m$. Es folgt

$$\frac{|F(x_1^{(8)})|}{m} \leq \frac{5}{5,8} 10^{-8} < 1 \cdot 10^{-8} = c.$$

Damit haben wir das Ergebnis

$$\xi = -0,68096220 \pm 10^{-8}.$$

Für die Größen $\frac{F^{(l)}(x)}{l!}$ ergeben sich in dem Bereich \mathfrak{B} die oberen Schranken $\frac{M_2}{2!} = 3$, $\frac{M_3}{3!} = 12$, $\frac{M_4}{4!} = 16$, $\frac{M_5}{5!} = 8$ ($F^{(l)}(x) \leq M_l$, $x \in \mathfrak{B}$). Damit erhält man (mit $c = 10^{-8}$) nach leichter Rechnung aus (5.16) und (5.17)

die Konstanten ¹	und	die Fehlerabschätzungen
$K_2 = 0,5173$		$ x_n - \xi \leq 1,93 [0,52 \cdot 10^{-8}]^{2^n}$
$K_3 = 2,61$		$ x_n - \xi \leq 0,62 [1,62 \cdot 10^{-8}]^{3^n}$
$K_4 = 23,1$		$ x_n - \xi \leq 0,36 [2,85 \cdot 10^{-8}]^{4^n}$
$K_5 = 493,9$		$ x_n - \xi \leq 0,22 [4,72 \cdot 10^{-8}]^{5^n}$

für die Iterationsverfahren 2. bis 5. Ordnung, wobei $x_0 = -0,68096220$ gesetzt ist.

¹ Die verhältnismäßig großen Beträge der Ableitungen $F^{(l)}(x)$ bei diesem Beispiel wirken sich auf die Konvergenz und die Abschätzungen bei höheren Verfahren ungünstig aus.

Tabelle 2. Schema zur Bestimmung einer Nullstelle der Gleichung

$$x^7 + 5x^6 + 3x^5 + 2x^4 + 4x^3 + 2x^2 + 6x + 4 = 0$$

durch ein Iterationsverfahren k -ter Ordnung: $x_1 = x_0 + v_k$, $x_0 = -0,75$

1	5	3	2	4	2	6	4
	-0,75	-3,1875	0,140625	-1,6054688	-1,7958985	-0,1530761	-4,3851929
1	4,25	-0,1875	2,140625	2,3945312	0,2041015	5,8469239	-0,3851929
	-0,75	-2,625	2,109375	-3,1875000	0,5947266	-0,5991211	$v_2 = 0,0734008$
1	3,50	-2,8125	4,250000	-0,7929688	0,7988281	(5,2478028)	
	-0,75	-2,0625	3,656250	-5,9296875	5,0419922		0,0314685
1	2,75	-4,8750	7,90625	-6,7226563	5,8408203	0,4287209	-0,3537244
	-0,75	-1,5	4,78125	-9,515625	-1,0948300		0,0215626
1	2	-6,375	12,68750	-16,2382813	4,7459903	0,3199002	$v_3 = 0,0674043$
	-0,75	-0,9375	5,484375	1,2591637	-1,0379315		$v_4 = 0,0692919$
1	1,25	-7,3125	18,1718750	-14,9791176	4,8028888	0,3328013	0,0230604
	-0,75	-0,375	-0,5304875	1,2173704	-1,0365405		-0,3621325
1	0,5	-7,6875	17,6413875	-15,0209109	4,8042798	0,3315265	0,0228775
	-0,75	-0,0172604	-0,5319474	1,2178853	-1,0370292		-0,3623154
1	-0,25	-7,7047604	17,6399276	-15,0203960	4,8037911	0,3316605	0,0228983
	0,0690374	-0,0124932	-0,5315875	1,2178396	-1,0369722		-0,3622946
1	-0,1809626	-7,6999932	17,6402875	-15,0204417	4,8038481	0,3316452	0,0228959
	0,0690378	-0,0124932	-0,5315906	1,2178464	-1,0369778		-0,3622970
1	-0,1809622	-7,6999932	17,6402844	-15,0204349	4,8038425	0,3316467	0,0228962
							$v_5 = 0,0690378 = v_4$

wie Verfahren 9. Ordnung

Literatur

- [1] SCHRÖDER, E.: Über unendlich viele Algorithmen zur Auflösung der Gleichungen. Math. Ann. **2**, 317—365 (1870).
- [2] BODEWIG, E.: Konvergenztypen und das Verhalten von Approximationen in der Nähe einer mehrfachen Wurzel einer Gleichung. Z. angew. Math. Mech. **29**, 44—51 (1949).
- [3] LUDWIG, R.: Über Iterationsverfahren für Gleichungen und Gleichungssysteme. Z. angew. Math. Mech. **34**, 210—225 (1954).
- [4] COLLATZ, L.: Näherungsverfahren höherer Ordnung für Gleichungen in Banach-Räumen. Arch. Rational Mech. Anal. **2**, 66—75 (1958).
- [5] ZURMÜHL, R.: Praktische Mathematik. Berlin-Göttingen-Heidelberg: Springer 1957.
- [6] HARTREE, D. R.: Notes on iterative processes. Proc. Cambridge phil. Soc. **45**, 230—236 (1949).
- [7] DOMB, C.: On iterative solutions of algebraic equations. Proc. Cambridge phil. Soc. **45**, 237—240 (1949).
- [8] KISS, I.: Über eine Verallgemeinerung des Newtonschen Näherungsverfahrens. Z. angew. Math. Mech. **34**, 68—69 (1954).
- [9] ZAJTA, A.: Untersuchungen über die Verallgemeinerung der Newton-Raphson-schen Wurzelapproximation. Acta techn. acad. scient. Hungaricae **15**, 233—260 (1956).

Mathematisches Institut A
der Technischen Hochschule
Stuttgart

(Eingegangen am 17. Juli 1959)

A Characterization of Hypergeometric Functions

N. CHAKO & J. MEIXNER

Summary

A theory of special functions should derive as many of the numerous formulas and properties as possible from some guiding principles. One of these principles is the F -equation, which TRUESDELL has put forward, and which is satisfied by all simple special functions. Many results about addition theorems, generating functions, integral representations *etc.* follow immediately from the F -equation. One disadvantage, however, is that the F -equation does not constitute a specific property of the simple special functions, generating instead functions of a much more general character.

By assuming two functional relations one of which can be transformed into an ascending and the other into a descending F -equation, one obtains as solutions of these functional relations the hypergeometric functions with all their special and confluent cases, and thus one has a new characterization of the hypergeometric functions.

1. Introduction

By the simple special functions of mathematical physics we understand the hypergeometric functions with their special and limiting cases, examples being Legendre functions, Bessel functions and confluent hypergeometric functions.

A common property of these simple special functions is that if they are written in appropriate variables, they become solutions of linear homogeneous differential equations of second order with polynomial coefficients. But this is a property common also to other special functions, such as Mathieu and spheroidal functions.

Another common property is the existence of linear three-term recurrence relations between three contiguous functions, where contiguous functions are defined as those differing by one in some parameter. The coefficients in these recurrence relations are polynomials in the independent variable and in the parameters, if both are suitably defined. But without specifying further the nature of the coefficients, these recurrence relations do not characterize of the simple special functions; such relations also hold true, for instance, for the Fourier expansion coefficients of Mathieu functions.

TRUESDELL* has drawn attention to the fact that the simple special functions, after appropriate transformation, satisfy a functional equation of the form

$$\frac{dF(z, \alpha)}{dz} = F(z, \alpha + 1) \quad (\alpha = \alpha_0, \alpha_0 \pm 1, \alpha_0 \pm 2, \dots), \quad (1)$$

* TRUESDELL, C.: An Essay toward a Unified Theory of Special Functions. Princeton University Press 1948.

which is called TRUESDELL's F -equation. This property (1) is remarkable and powerful because it permits one to derive in a brief and straightforward fashion many theorems for simple special functions; in particular, addition theorems, generating functions, integral relations. But this F -equation is not a characterizing property of the simple special functions either, since the generalized zeta function $\zeta(\alpha, z)$, after multiplication by $e^{i\alpha\pi}\Gamma(\alpha)$, satisfies (1) without belonging to the functions of hypergeometric type.

However, the power of functional equations like (1) can be applied to a characterization of the simple special functions, if one recalls that the simple special functions always satisfy two simultaneous functional equations of the types

$$\frac{d\gamma(z, \alpha)}{dz} = U_1(z, \alpha) \gamma(z, \alpha) + V_1(z, \alpha) \gamma(z, \alpha + 1), \quad (2)$$

$$\frac{d\gamma(z, \alpha)}{dz} = U_2(z, \alpha) \gamma(z, \alpha) + V_2(z, \alpha) \gamma(z, \alpha - 1). \quad (3)$$

By themselves, these equations do not give a sufficient restriction on the possible solutions of (2) and (3), so long as we do not stipulate additional conditions. We require therefore, that there exist functions $z(\xi_i)$ and $\Phi_i(z(\xi_i), \alpha)$, $i=1, 2$, such that the functions $Y_i(\xi_i, \alpha) = \Phi_i(z(\xi_i), \alpha) \gamma(z(\xi_i), \alpha)$ satisfy the ascending F -equation

$$\frac{dY_1(\xi_1, \alpha)}{d\xi_1} = Y_1(\xi_1, \alpha + 1)$$

and the descending F -equation

$$\frac{dY_2(\xi_2, \alpha)}{d\xi_2} = Y_2(\xi_2, \alpha - 1).$$

Even so, too much arbitrariness remains, and therefore we require in addition that there be two linearly independent solutions of (2) and (3) with the above mentioned properties.

There is no loss in generality if we assume $U_1(z, \alpha)=0$, $V_1(z, \alpha)=1$, which simply means that we transform the equations (2) and (3) by taking $Y_1(\xi_1, \alpha)$ as new dependent variable.

Changing the notation slightly, we therefore seek functions $\gamma(z, \alpha)$, $\alpha=\alpha_0$, $\alpha_0 \pm 1$, $\alpha_0 \pm 2, \dots$, with the following properties*:

1. $\gamma(z, \alpha)$ is a solution of the ascending F -equation

$$\frac{d\gamma(z, \alpha)}{dz} = \gamma(z, \alpha + 1). \quad (4)$$

2. There exists a function $\Phi(z, \alpha)$ and a variable $\xi=\xi(z)$ with $d\xi/dz \neq 0$ in some domain of the complex variable z such that

$$\gamma(\xi, \alpha) = \Phi(z(\xi), \alpha) \gamma(z(\xi), \alpha) \quad (5)$$

is a solution of the descending F -equation

$$\frac{dY(\xi, \alpha)}{d\xi} = Y(\xi, \alpha - 1). \quad (6)$$

3. There exists a second function $\gamma_1(z, \alpha)$, which is linearly independent of $\gamma(z, \alpha)$ and has the same properties 1. and 2.

* This is a generalization of a problem which TRUESDELL indicated in Theorem 6.1 of his monograph.

2. Derivation of the general solution

We note first the following

Lemma. *If $y(z, \alpha)$ is a solution of the ascending F-equation, so also is*

$$k^\alpha y(kz + l, \alpha + a).$$

If $Y(\xi, \alpha)$ is a solution of the descending F-equation, so also is

$$k^{-\alpha} Y(k\xi + l, \alpha + a).$$

Here k, l , and a are arbitrary constants with the restriction $k \neq 0$.

The proof of this lemma is elementary.

In the following considerations, for simplicity of notation, we often omit the variable z from the arguments of the functions $y(z, \alpha)$ and $\Phi(z, \alpha)$. A prime means differentiation with respect to the variable z .

By inserting (5) into (6), we obtain the equation

$$[\Phi'(\alpha) y(\alpha) + \Phi(\alpha) y'(\alpha)] \frac{dz}{d\xi} = \Phi(\alpha - 1) y(\alpha - 1). \quad (7)$$

The equations (7) and (4) are not compatible unless some conditions on the functions $\Phi(\alpha)$ and $z(\xi)$ hold. To derive them, we replace α by $\alpha - 1$ in (4) and by $\alpha + 1$ in (7). Let us denote the resulting equations by (4') and (7'). Then we eliminate $y(\alpha - 1)$ from (4') and (7) and $y(\alpha + 1)$ from (4) and (7'). This elimination results in two linear and homogeneous differential equations of second order for $y(z, \alpha)$. They are

$$\begin{aligned} & [\Phi''(\alpha) y(\alpha) + 2\Phi'(\alpha) y'(\alpha) + \Phi(\alpha) y''(\alpha)] \frac{dz}{d\xi} + \\ & + [\Phi'(\alpha) y(\alpha) + \Phi(\alpha) y'(\alpha)] \left[\frac{d}{dz} \left(\frac{dz}{d\xi} \right) - \frac{\Phi'(\alpha - 1)}{\Phi(\alpha - 1)} \frac{dz}{d\xi} \right] = \Phi(\alpha - 1) y(\alpha), \end{aligned} \quad (8)$$

$$[\Phi(\alpha + 1) y''(\alpha) + \Phi'(\alpha + 1) y'(\alpha)] \frac{dz}{d\xi} - \Phi(\alpha) y(\alpha) = 0. \quad (9)$$

Since both differential equations are satisfied by the same function $y(z, \alpha)$ and by the same linearly independent function $y_1(z, \alpha)$, they must be identical. By comparing their coefficients, one therefore obtains the following two equations for $\Phi(\alpha)$ and $z(\xi)$:

$$\frac{\Phi'(\alpha + 1)}{\Phi(\alpha + 1)} - 2 \frac{\Phi'(\alpha)}{\Phi(\alpha)} + \frac{\Phi'(\alpha - 1)}{\Phi(\alpha - 1)} = \frac{(dz/d\xi)'}{dz/d\xi}, \quad (10)$$

$$\frac{\Phi''(\alpha)}{\Phi(\alpha)} \frac{dz}{d\xi} + \frac{\Phi'(\alpha)}{\Phi(\alpha)} \frac{d}{dz} \left(\frac{dz}{d\xi} \right) - \frac{\Phi'(\alpha) \Phi'(\alpha - 1)}{\Phi(\alpha) \Phi(\alpha - 1)} \frac{dz}{d\xi} + \frac{\Phi(\alpha)}{\Phi(\alpha + 1)} - \frac{\Phi(\alpha - 1)}{\Phi(\alpha)} = 0. \quad (11)$$

The first equation constitutes a linear difference equation in α of second order for the function $\Phi'(\alpha)/\Phi(\alpha)$, which is solved by a polynomial in α of second degree,

$$\frac{\Phi'(\alpha)}{\Phi(\alpha)} = p'_0(z) + \alpha p'_1(z) + \alpha^2 p'_2(z), \quad (12)$$

with coefficients $p'_0(z), p'_1(z), p'_2(z)$ which may depend on z . Now the integration with respect to z can be performed. It gives

$$\Phi(\alpha) = q(\alpha) \exp [p_0(z) + \alpha p_1(z) + \alpha^2 p_2(z)]. \quad (13)$$

The factor $q(\alpha)$ is the integration constant and thus depends only on α .

Due to (12), the first three terms in (11) form a polynomial of at most fourth degree in α . Considering (11) as a difference equation for $\Phi(\alpha)/\Phi(\alpha+1)$, one observes that this expression has to be a polynomial of fifth degree or less in α . On the other hand, one obtains from (13)

$$\frac{\Phi(\alpha)}{\Phi(\alpha+1)} = \frac{q(\alpha)}{q(\alpha+1)} \exp[-p_1(z) - (2\alpha+1)p_2(z)], \quad (14)$$

which is a polynomial in α for all z , only if $p_2(z)$ is a constant. Without loss of generality, one can put this constant equal to zero by redefining $q(\alpha)$ as the $q(\alpha)$ in (13) multiplied by $\exp(\alpha^2 p_2)$.

Thus, putting $p_2(z)=0$, one concludes from (12) that $\Phi'(\alpha)/\Phi(\alpha)$ is a polynomial of degree at most one in α , and therefore the right member of (10) vanishes. This is equivalent to z 's being a linear function of ξ and vice versa. By the lemma, one can restrict consideration to the special case $z=\xi$ and then arrive at the general function $Y(\xi, \alpha)$ connected with $y(z, \alpha)$ by the transformation given in the lemma.

What remains now is the discussion of equation (11). Making use of (13), of $p_2=0$ and $z=\xi$, we simplify (11) to the form

$$\begin{aligned} p_0''(z) + p_0'(z)p_1'(z) + \alpha[p_1''(z) + p_1'(z)^2] + \\ + \left[\frac{q(\alpha)}{q(\alpha+1)} - \frac{q(\alpha-1)}{q(\alpha)} \right] \exp(-p_1(z)) = 0. \end{aligned} \quad (15)$$

One concludes from this equation that $q(\alpha)/q(\alpha+1)$ is a polynomial of at most the second degree in α , namely

$$-\frac{q(\alpha)}{q(\alpha+1)} = q_0 + s_1\alpha + r_2\alpha(\alpha+1) \quad (16)$$

with constants coefficients q_0, s_1, r_2 . Inserting (16) into (15) then gives

$$p_0''(z) + p_0'(z)p_1'(z) = s_1 \exp[-p_1(z)], \quad (17)$$

$$p_1''(z) + p_1'(z)^2 = 2r_2 \exp[-p_1(z)]. \quad (18)$$

The solutions of these differential equations are

$$p_1 = \ln(r_0 + r_1 z + r_2 z^2), \quad p_0' = \frac{s_0 + s_1 z}{r_0 + r_1 z + r_2 z^2} \quad (19)$$

with new constants r_0, r_1, s_0 .

By use of (12), (14), (16), (19), the differential equation (9) becomes

$$\begin{aligned} (r_0 + r_1 z + r_2 z^2) y''(z, \alpha) + [s_0 + s_1 z + (\alpha+1)(r_1 + 2r_2 z)] y'(z, \alpha) + \\ + [q_0 + s_1 \alpha + r_2 \alpha(\alpha+1)] y(z, \alpha) = 0 \end{aligned} \quad (20)$$

with arbitrary constants $r_0, r_1, r_2, s_0, s_1, q_0$. This is a differential equation of hypergeometric type. The same differential equation may be obtained from (8).

3. Discussion of the solution

While it is clear from (20) that the functions which satisfy the postulates mentioned in Section 2 are hypergeometric functions, including their special and limiting cases, it is worthwhile to give explicit expressions for $y(z, \alpha)$, $\Phi(z, \alpha)$, $Y(z, \alpha)$ in familiar notations, because, by TRUESEDELL'S results, to exhibit such expressions is tantamount to giving explicit generating functions and addition theorems for the various simple special functions.

We must distinguish several cases according to the character of the polynomial coefficient $r_0 + r_1 z + r_2 z^2$. It can have two different roots, two equal roots, one single root ($r_2 = 0$), or none at all ($r_2 = r_1 = 0$, $r_0 \neq 0$).

In each case the procedure is to give a suitable solution of (20) in terms of the conventional notations of the simple special functions, to multiply it by a factor which depends only on α such that (4) is satisfied — this is then the $y(z, \alpha)$ — to evaluate $q(\alpha)$ from (16) and $\Phi(z, \alpha)$ from (13), and finally to write down the $Y(z, \alpha)$ according to (5).

Case 1. The polynomial $r_0 + r_1 z + r_2 z^2$ has two distinct roots, which can be assumed to be at $z=0$ and $z=1$ (see the Lemma); that is, $r_0=0$, $r_1=-r_2$. Furthermore, one can assume $r_2=1$. For convenience s_0, s_1, q_0 are expressed by other constants a, b, c in the following way:

$$s_1 = a + b - 1, \quad a_0 = a b, \quad s_0 = 1 - c.$$

Then the differential equation (20) reads

$$z(z-1)y''(z, \alpha) + [(a+b+2\alpha+1)z - (c+\alpha)]y'(z, \alpha) + (a+\alpha)(b+\alpha)y(z, \alpha) = 0. \quad (21)$$

This is the hypergeometric differential equation in normal form with the parameters $a+\alpha, b+\alpha, c+\alpha$.

From (16) one obtains

$$q(\alpha) = \frac{(-1)^\alpha}{\Gamma(a+\alpha)\Gamma(b+\alpha)}, \quad (22)$$

apart from a trivial constant factor, which can be omitted. (19) gives

$$\begin{aligned} p_1(z) &= \ln[z(z-1)], \\ \exp p_0(z) &= z^{c-1}(z-1)^{a+b-c}, \end{aligned} \quad (23)$$

again omitting a trivial constant factor. Inserting these results into (13) leads to

$$\Phi(\alpha) = \frac{1}{\Gamma(a+\alpha)\Gamma(b+\alpha)} z^{c+\alpha-1} (1-z)^{a+b+\alpha-c}, \quad (24)$$

up to a constant factor.

Not every solution of the differential equation (21) is a solution of the ascending F -equation (4). But from well known properties of the hypergeometric function, one verifies immediately that

$$y(z, \alpha) = \frac{\Gamma(a+\alpha)\Gamma(b+\alpha)}{\Gamma(c+\alpha)} {}_2F_1(a+\alpha, b+\alpha; c+\alpha; z) \quad (25)$$

has all the required properties.

Finally one obtains from (5) with $z=\xi$

$$Y(z, \alpha) = \frac{1}{\Gamma(c+\alpha)} z^{c+\alpha-1} (1-z)^{a+b+\alpha-c} {}_2F_1(a+\alpha, b+\alpha; c+\alpha; z). \quad (26)$$

A whole set of other solutions is obtained if one replaces the hypergeometric function $F(a+\alpha, b+\alpha; c+\alpha; z)$ by one of the functions

$$\begin{aligned} & (1-z)^{-a-\alpha} {}_2F_1\left(a+\alpha, c-b; c+\alpha; \frac{z}{z-1}\right), \\ & (1-z)^{c-a-b-\alpha} {}_2F_1(c-a, c-b; c+\alpha; z), \\ & \Gamma(c+\alpha) \Gamma(c-a-b-\alpha) {}_2F_1(a+\alpha, b+\alpha; a+b+\alpha-c+1; 1-z), \\ & \frac{\Gamma(c+\alpha)}{\Gamma(b+\alpha)} (-1)^\alpha z^{-a-\alpha} {}_2F_1(a+\alpha; 1-c+a; 1-b+a; z^{-1}). \end{aligned}$$

But, of course, all these solutions can be expressed in terms of any two which are linearly independent.

In each of these examples, the hypergeometric function ${}_2F_1(A, B; C; u)$ can be replaced by

$${}_2\Psi_1(A, B; C; u) = \frac{\Gamma(A') \Gamma(B') \Gamma(C-1)}{\Gamma(A) \Gamma(B) \Gamma(C'-1)} u^{C'-1} {}_2F_1(A', B'; C'; u)$$

with $A'=A-C+1$, $B'=B-C+1$; $C'=2-C$. Thus one obtains other forms of the solutions.

Case 2. The polynomial $r_0+r_1z+r_2z^2$ has two equal roots. Using the lemma, one can assume that $r_0=r_1=0$, $r_2=1$. We put $q_0=a(a-c+1)$, $s_1=2a-c$. Then the following results, which can be derived similarly as in Case 1, hold:

$$\begin{aligned} & z^2 y''(z, \alpha) + [s_0 + (2a + 2\alpha - c + 2)z] y'(z, \alpha) + \\ & + (a + \alpha)(a + \alpha - c + 1) y(z, \alpha) = 0, \\ & y(z, \alpha) = (-1)^\alpha z^{-a-\alpha} \Gamma(a+\alpha) {}_1F_1\left(a+\alpha; c; \frac{s_0}{z}\right), \\ & \Phi(z, \alpha) = \frac{(-1)^\alpha}{\Gamma(a+\alpha) \Gamma(a-c+\alpha+1)} e^{-\frac{s_0}{z}} z^{2a+2\alpha-c}, \\ & y(z, \alpha) = \frac{z^{a+\alpha-c}}{\Gamma(a+\alpha-c+1)} e^{-\frac{s_0}{z}} {}_1F_1\left(a+\alpha; c; \frac{s_0}{z}\right). \end{aligned}$$

Instead of the confluent hypergeometric function, ${}_1F_1\left(a+\alpha; c; \frac{s_0}{z}\right)$, one can also take the function

$${}_1\Psi_1\left(a+\alpha; c; \frac{s_0}{z}\right) = \frac{\Gamma(a+\alpha-c+1) \Gamma(c-1)}{\Gamma(a+\alpha) \Gamma(1-c)} \left(\frac{z}{s_0}\right)^{c-1} {}_1F_1\left(a+\alpha-c+1; 2-c; \frac{s_0}{z}\right).$$

A simple special case is $s_0=0$, which however, will not be discussed further.

Case 3. The polynomial $r_0+r_1z+r_2z^2$ has only one single root. One therefore assumes $r_2=0$, $r_1=1$, $r_0=0$. Here we must distinguish between $s_1 \neq 0$ and $s_1=0$.

Subcase 3a. $s_1 \neq 0$. Then one can put $s_1 = -1$ without loss of generality. Introducing the constants $c = s_0 + 1$, $a = -q_0$, we obtain

$$z y''(z, \alpha) + (c + \alpha - z) y'(z, \alpha) - (a + \alpha) y(z, \alpha) = 0,$$

$$y(z, \alpha) = \frac{\Gamma(a + \alpha)}{\Gamma(c + \alpha)} {}_1F_1(a + \alpha; c + \alpha; z),$$

$$\Phi(z, \alpha) = \frac{1}{\Gamma(a + \alpha)} e^{-z} z^{c + \alpha - 1},$$

$$y(z, \alpha) = \frac{1}{\Gamma(c + \alpha)} e^{-z} z^{c + \alpha - 1} {}_1F_1(a + \alpha; c + \alpha; z).$$

A further solution is obtained with the same replacement of the confluent hypergeometric function ${}_1F_1$ by ${}_1\Psi_1$, as given in Case 2.

Subcase 3b. Let $s_1 = 0$. Then

$$z y''(z, \alpha) + (s_0 + \alpha + 1) y'(z, \alpha) + q_0 y(z, \alpha) = 0,$$

$$y(z, \alpha) = (-1)^\alpha q_0^{\alpha/2} z^{-(s_0 + \alpha)/2} J_{s_0 + \alpha}(\sqrt{4q_0 z}),$$

$$\Phi(z, \alpha) = (-q_0)^{-\alpha} z^{s_0 + \alpha},$$

$$Y(z, \alpha) = q_0^{-\alpha/2} z^{(s_0 + \alpha)/2} J_{s_0 + \alpha}(\sqrt{4q_0 z}).$$

The Bessel function in these formulas can be replaced by the Neumann function or by one of the Hankel functions with the same argument and order.

Case 4. The polynomial $r_0 + r_1 z + r_2 z^2$ reduces to a constant, which is assumed to be 1. Then we again have to distinguish between the cases $s_1 \neq 0$ and $s_1 = 0$.

Subcase 4a. Let $s_1 \neq 0$. Without loss of generality we assume $s_1 = 1$, $s_0 = 0$ and put $q_0 = 1 + a$. Then

$$y''(z, \alpha) + z y'(z, \alpha) + (1 + a + \alpha) y(z, \alpha) = 0,$$

$$y(z, \alpha) = (-1)^\alpha \exp\left(-\frac{z^2}{4}\right) D_{a + \alpha}(z),$$

$$\Phi(z, \alpha) = \frac{(-1)^\alpha}{\Gamma(a + \alpha + 1)} \exp\left(\frac{z^2}{2}\right),$$

$$Y(z, \alpha) = \frac{1}{\Gamma(a + \alpha + 1)} \exp\left(\frac{z^2}{4}\right) D_{a + \alpha}(z).$$

Instead of the function of the parabolic cylinder $D_{a + \alpha}(z)$, one can also take the following function:

$$\Gamma(a + \alpha + 1) i^\alpha D_{-a - \alpha - 1}(iz).$$

Subcase 4b. Let $s_1 = 0$. In this case the differential equation (20) become independent of α . Considering all possible values of s_0 and q_0 , one arrives at the following pair of solutions ($a \neq 0$):

$$Y(z, \alpha) = a^\alpha e^{az}, \quad \tilde{Y}(z, \alpha) = a^\alpha (az + \alpha) e^{az},$$

$$\Phi(z, \alpha) = (-1)^\alpha a^{-2\alpha} e^{-2az},$$

$$y(z, \alpha) = (-1)^\alpha a^{-\alpha} e^{-az}, \quad \tilde{y}(z, \alpha) = (-1)^\alpha a^{-\alpha} (az + \alpha) e^{-az}.$$

4. Remarks and discussion

In the foregoing results a constant factor (constant with respect to both variables z and α) can always be added in $y(z, \alpha)$ as well as in $Y(z, \alpha)$. Furthermore, the results remain correct if $y(z, \alpha)$ and $Y(z, \alpha)$ are multiplied by a periodic function of α with period 1. One also sees immediately that $y(z, \alpha)$, $\Phi(z, \alpha)$, $Y(z, \alpha)$ can in turn be changed into $Y(z, -\alpha)$, $\Phi(z, -\alpha)^{-1}$, $y(z, -\alpha)$ to give a solution with the properties mentioned in Section 1. Although one does not obtain new solutions in this way, the totality of solutions having been exhausted by the cases 1 to 4 to within linear transformations of the independent variable, it may, however, be convenient to put the solutions into a different form in this way.

Special situations, such as terminating sequences $y(z, \alpha)$, are not explicitly treated; it is easy to make the necessary modifications in every case.

Of course, the limiting cases 2 to 4 can also be obtained by going to the respective limits not in the differential equations, but in the solutions themselves, $y(z, \alpha)$, $\Phi(z, \alpha)$, $Y(z, \alpha)$, of Case 1. Therefore everything is theoretically contained in Case 1, explicitly in the formulas (24), (25), (26) and in those which arise by replacing α by $-\alpha$ and y, Φ, Y by Y, Φ^{-1}, y .

We repeat that the solutions $y(z, \alpha)$ of the ascending F -equation written down explicitly in the Cases 1 to 4 satisfy the functional equation (7) with $dz/d\xi = 1$, and furthermore we remark that (7) and (4) can be combined into the difference equation or the recurrence relation

$$y(z, \alpha + 1) + \frac{\Phi'(z, \alpha)}{\Phi(z, \alpha)} y(z, \alpha) - \frac{\Phi(z, \alpha - 1)}{\Phi(z, \alpha)} y(z, \alpha - 1) = 0.$$

Legendre functions are special cases of hypergeometric functions and can be expressed in terms of these in various ways. This gives rise to many functions $y(z, \alpha)$ that can be expressed by Legendre functions. However, we shall not go into the details here.

Queen's College, Queens, N.Y.
and

Technische Hochschule Aachen

(Received June 26, 1959)

EDITORIAL BOARD

R. BERKER
Technical University
Istanbul

L. CESARI
Research Institute for Advanced Study
Baltimore, Maryland

L. COLLATZ
Institut für Angewandte Mathematik
Universität Hamburg

A. ERDÉLYI
California Institute of Technology
Pasadena, California

J. L. ERICKSEN
The Johns Hopkins University
Baltimore, Maryland

G. FICHERA
Mathematics Research Center
U. S. Army
University of Wisconsin
Madison, Wisconsin

R. FINN
Stanford University
California

HILDA GEIRINGER
Harvard University
Cambridge, Massachusetts

H. GÖRTLER
Institut für Angewandte Mathematik
Universität Freiburg i. Br.

D. GRAFFI
Istituto Matematico „Salvatore Pincherle“
Università di Bologna

A. E. GREEN
King's College
Newcastle-upon-Tyne

J. HADAMARD
Institut de France
Paris

L. HÖRMANDER
Department of Mathematics
University of Stockholm

M. KAC
Cornell University
Ithaca, New York

E. LEIMANIS
University of British Columbia
Vancouver

A. LICHNEROWICZ
Collège de France
Paris

C. C. LIN
Massachusetts Institute of Technology
Cambridge, Massachusetts

W. MAGNUS
Institute of Mathematical Sciences
New York University
New York City

G. C. McVITTIE
University of Illinois Observatory
Urbana, Illinois

J. MEIXNER
Institut für Theoretische Physik
Technische Hochschule Aachen

C. MIRANDA
Istituto di Matematica
Università di Napoli

C. B. MORREY
University of California
Berkeley, California

C. MÜLLER
Mathematisches Institut
Technische Hochschule Aachen

W. NOLL
Carnegie Institute of Technology
Pittsburgh, Pennsylvania

A. OSTROWSKI
Mathematics Research Center
U. S. Army
University of Wisconsin
Madison, Wisconsin

R. S. RIVLIN
Division of Applied Mathematics
Brown University
Providence, Rhode Island

M. M. SCHIFFER
Stanford University
California

J. SERRIN
Institute of Technology
University of Minnesota
Minneapolis, Minnesota

E. STERNBERG
Division of Applied Mathematics
Brown University
Providence, Rhode Island

R. TIMMAN
Instituut voor Toegepaste Wiskunde
Technische Hogeschool, Delft

R. A. TOUPIN
Naval Research Laboratory
Washington 25, D.C.

C. TRUESDELL
801 North College Avenue
Bloomington, Indiana

H. VILLAT
47, bd. A. Blanqui
Paris XIII

CONTENTS

TRUESDELL, C., Invariant and Complete Stress Functions for General Continua	1
GENENSKY, S. M., & R. S. RIVLIN, Infinitesimal Plane Strain in a Network of Elastic Cords	30
EHRMANN, H., Iterationsverfahren mit veränderlichen Operatoren . .	45
EHRMANN, H., Konstruktion und Durchführung von Iterationsverfahren höherer Ordnung	65
CHAKO, N., & J. MEIXNER, A Characterization of Hypergeometric Functions	89